AUTOMATING THE ANALYSIS AND IMPROVEMENT OF DYNAMIC PROGRAMMING ALGORITHMS WITH APPLICATIONS TO NATURAL LANGUAGE PROCESSING

by

Tim Vieira

A dissertation submitted to The Johns Hopkins University in conformity with the requirements for the degree of Doctor of Philosophy.

Baltimore, Maryland

July, 2023

© 2023 Tim Vieira

All rights reserved

Abstract

This thesis develops a system for automatically analyzing and improving dynamic programs, such as those that have driven progress in natural language processing and computer science, more generally, for decades. Finding a correct program with the optimal asymptotic runtime can be unintuitive, time-consuming, and error-prone. This thesis aims to automate this laborious process. To this end, we develop an approach based on

- a high-level, domain-specific language called Dyna for concisely specifying dynamic programs
- 2. a general-purpose solver to efficiently execute these programs
- 3. a static analysis system that provides type analysis and worst-case time/s-pace complexity analyses
- 4. a rich collection of meaning-preserving transformations to programs, which systematizes the repeated insights of numerous authors when speeding up algorithms in the literature
- 5. a search algorithm for identifying a good sequence of transformations that

reduce the runtime complexity given an initial, correct program

We show that, in practice, automated search—like the mental search performed by human programmers—can find substantial improvements to the initial program. Empirically, we show that many speed-ups described in the NLP literature could have been discovered automatically by our system. We provide a freely available prototype system at https://github.com/timvieira/dyna-pi.

Primary Reader and Advisor:

Jason Eisner (Johns Hopkins University)

Secondary Readers:

David Chiang (University of Notre Dame)

Scott Smith (Johns Hopkins University)

Acknowledgments

Academic

I am incredibly grateful to my advisor, Jason Eisner, for taking me on as a student and patiently working with me all these years. I have learned *so much* about how to teach, present, write, think, get the details right—and have fun doing it! My PhD has been a truly transformational experience thanks to Jason.

I want to thank Ryan Cotterell, who acted as an unofficial secondary advisor in my last two years. Without Ryan's support, knowledge, patience, cheerleading, and insights, I would not have been able to write this dissertation.

I want to thank Wes Filardo, who taught me an incredible amount of what underlies this thesis. Sitting next to such a knowledgeable and generous human made my PhD much more valuable and enjoyable.

I want to thank my committee members—Jason Eisner, David Chiang, and Scott Smith—for patiently reading this *long* manuscript and providing useful and insightful feedback. Their feedback has greatly improved this manuscript. Any mistakes are my own.

I want to thank Andrew McCallum and Dan Roth for employing me as a research programmer, which set me on track to do a PhD. Andrew even provided the crucial introduction to Jason Eisner!

I want to thank Gerald DeJong, Jeff Erickson, Viraj Kumar, and Cinda Heeren for their support during my undergraduate years. Viraj is the reason why I first felt empowered to go to graduate school—I am very grateful to him.

My PhD experience would not have been the same without the extraordinary academic community I have been fortunate to be part of.

- Collaborators: Afra Amini, Alexandra Butoi, Andreas Opedal, Anej Svete,
 Benjamin Dayan, Brian DuSell, Clara Meister, Clemente Pasti, He He,
 Jiarong Jiang, Jiawei Huang, Matthew Francis-Landau, Qi He, Ran Zmigrod, Tiago Pimentel, and Wes Filardo.
- Labmates (JHU): Adam Teichert, Brian Lu, Chu-Cheng Lin, Darcey Riley, Dingquan Wang, Hongyuan Mei, Juneki Hong, Leo Du, Matthew Francis-Landau, Matt Gormley, Michael Paul, Nick Andrews, Ryan Cotterell, Sabrina Mielke, and Wes Filardo.
- Labmates (UMass): Greg Druck, Sameer Singh, Sebastian Riedel, Michael
 Wick, Kedar Bellare, Adam Saunders, Rachel Shorey, David Mimno, Pallika Kanani, Anton Bakalov, Karl Schultz, Limin Yao, Meagan Day, Laura
 Dietz, Hanna Wallach, and Jason Naradowsky.
- Labmates (UIUC): Alex Klementiev, Dan Goldwasser, Jeff Pasternack, Kevin Small, Lev Ratinov, Mark Sammons, Michael Connor, Ming-Wei

Chang, Nick Rizzolo, Quang Do, Rajhans Samdani, Subhro Roy, V.G.Vinod Vydiswaran, Vivek Srikumar, and Yuancheng Tu.

- Broader JHU: Ben Van Durme, Colin Lea, Elizabeth Salesky, Frank Ferraro, Hainan Xu, Huda Khayrallah, Juri Ganitkevitch, Kim Franklin, Keith Levin, Sanjeev Khudanpur, Tongfei Chen, Mark Dredze, Naomi Saphra, Olivia Buzek, Pushpendre Rastogi, Rachel Rudinger, Ruth Scally, Sheng Zhang, Travis Wolfe, and Yanif Ahmad.
- Broader Community: Aaron Schein, David Bellanger, David Smith, Hal
 Daumé, Kianté Brantley, Mathieu Blondel, Nikos Karampatziakis, Phil
 Thomas, Pierre-Luc Bacon, Sasha Rush, and Vlad Niculae.

Non-Academic

I want to thank my parents, Marcia and Terry Vieira, for their unconditional love and support, for always prioritizing my education, and for their sacrifices to send me to private school and out-of-state college. I could not ask for better parents. Thank you, Mom and Dad!

I want to thank my incredible friends and family for all their love, support, and entertainment: Hanna Wallach, Maia Papaya, Marcia Vieira, Terry Vieira, Tamara Vieira, Roseli Fabbri, Rachael Wallach von Portheim, Luna Petunia, Robin Wallach, Rob Wallach, Parth Naik, Ryan Cotterell, Boris Umanov, Alex O'Neill, Keanu, Dan Muriello, Jason Naradowsky, Ian Baran, Dan Schreiber,

Abbie Jacobs, Aaron Schein, Clint Adams, Hal Daumé, and Amber Boydstun.

The biggest and most important thank you is to my partner, Hanna Wallach. She has been with me through my entire PhD journey—she even helped me with my application essay! I would not be here at the finish line without her love, encouragement, generosity, patience, and support. Thank you, Hanna.

Previous Publications

Portions of this thesis have been published in the following:

- Tim Vieira, Matthew Francis-Landau, et al. (2017). "Dyna: Toward a Self-Optimizing Declarative Language for Machine Learning Applications". In: Proceedings of the ACM SIGPLAN Workshop on Machine Learning and Programming Languages), pp. 8–17. ISBN: 978-1-4503-5071-6. URL: http://cs.jhu.edu/~jason/papers/#vieira-et-al-2017
- Tim Vieira, Ryan Cotterell, et al. (2021). "Searching for More Efficient Dynamic Programs". In: Findings of the Association for Computational Linguistics. URL: https://arxiv.org/abs/2109.06966

Contents

A l	bstra	ıct			ii
A	ckno	wledg	ments		v
P	revio	ous Pu	blications		ix
Li	ist of	Table	s	3	xvi
Li	ist of	Figur	res	X	vii
Li	ist of	Algor	ithms	XV	viii
1	Int	roduct	ion		1
2	A D	omain	a-Specific Language for Dynamic Programming		9
3	Not	tation,	Background, and Utilities		17
	3.1	Inter	nal Representation of Programs		18
		3.1.1	Notation and Representation		18
		3.1.2	Lookup and Update Represention		22
			3.1.2.1 Implicit Result Stream Data Structure		23

		3.1.2.2 Lookup $(\mathfrak{p}[x])$	24
		3.1.2.3 Update ($\mathfrak{p}[x] \oplus =$)	29
	3.2	General Mathematics	31
		3.2.1 Ordering Relations	32
		3.2.2 Functions	34
		3.2.3 Graphs	35
		3.2.4 Abstract Algebra	38
	3.3	Semirings: A Powerful Abstraction for Dynamic Programming	42
4	Son	nantics of Dyna Programs	54
•			
	4.1	Total Value of All Derivations	55
	4.2	Fixpoint Approximation	60
	4.3	Fixpoint Iteration	63
	4.4	Convergence Theory	65
5	Infe	erence Algorithms	71
	5.1	Algorithms for Fixpoint Iteration	73
		5.1.1 Naïve Fixpoint Iteration	73
		5.1.2 Semi-Naïve Fixpoint Iteration	75
		5.1.3 Prioritized Fixpoint Iteration	80
		5.1.4 Prioritization Functions	83
	5.2	Newton's Method	87
	5.3	Specialized Right-Linear Solver	92
	5.4	Hybrid Algorithm	96

6	Pro	ogram Analysis 10		
	6.1	Type A	Analysis	102
		6.1.1	Formulation	102
		6.1.2	Abstract Forward Chaining	106
			6.1.2.1 Type Inference Algorithm	107
			6.1.2.2 Relaxation	108
			6.1.2.3 Constraint Propagation	111
			6.1.2.4 Removing Redundant Simple Types	114
			6.1.2.5 Termination	115
	6.2	Type-	based Program Specialization	120
		6.2.1	Inferring a Type for Useful Items	122
		6.2.2	Specialization Transformation	123
		6.2.3	Abbreviation Transformation	125
	6.3	Auton	natic Space and Time Analysis	129
		6.3.1	Space Analysis	129
		6.3.2	Time Analysis	133
	6.4	Applio	cations	137
7	Pro	gram '	Transformations	142
	7.1	Inforn	nal Introduction	144
	7.2	Semai	ntics-Preserving Transformations	152
	7.3	Defini	tion-Introduction Transformation	156
	7.4		and Unfold Transformations	

		7.4.1	Basic Fol	d and Unfold	161
		7.4.2	Generaliz	zed Fold and Unfold	163
			7.4.2.1	Formal Definitions	166
			7.4.2.2	Challenges	168
			7.4.2.3	What Causes Partial Correctness?	169
			7.4.2.4	Practical Sufficient Conditions for Correctness	175
	7.5	Hoisti	ng Transf	formations	186
		7.5.1	Speculati	ion Transformation	186
		7.5.2	Left-Corr	ner Transformation	193
		7.5.3	Examples	s	196
	7.6	Discu	ssion		201
8	Sea	rching	g For Fast	ter Dynamic Programs	206
8				ter Dynamic Programs	
8			raph Sear	·	207
8		The G	raph Sear Prioritize	ch Problem	207 208
8	8.1	The G 8.1.1 8.1.2	raph Sear Prioritize Search S	cch Problem	207208210
8	8.1	The G 8.1.1 8.1.2	raph Sear Prioritize Search S Tunctions	ch Problem	207208210214
8	8.1	The G 8.1.1 8.1.2 Cost I	raph Sear Prioritize Search S Functions Fine-grai	cch Problem	207208210214217
8	8.1	The G 8.1.1 8.1.2 Cost I 8.2.1 8.2.2	Prioritize Search Sounctions Fine-grai	rch Problem	207208210214217219
9	8.1 8.2 8.3	The G 8.1.1 8.1.2 Cost I 8.2.1 8.2.2 Exper	Prioritize Search Sounctions Fine-grai	rch Problem	207208210214217219
	8.1 8.2 8.3 Rel a	The G 8.1.1 8.1.2 Cost I 8.2.1 8.2.2 Exper	raph Sear Prioritize Search S Functions Fine-grai Coarse-grai iments	rch Problem	207 208 210 214 217 219 222 227

10	10 Conclusion 23				
A	Infe	erence Algorithms	238		
В	Pro	gram Transformations	241		
	B.1	Strong Equivalence	242		
	B.2	Correctness of the Fold and Unfold Transformations	243		
	B.3	Speculation	264		
C	Pat	tern Matching Algorithms	266		
	C.1	Unification	266		
	C.2	Sub-Product Enumeration	268		
		C.2.1 Commutative Sub-Products	268		
		C.2.2 Noncommutative Sub-Products	269		
D	Dyn	na 2	27 0		
E	Pro	grams Used for Experiments	274		
	E.1	Chains (chain-05, chain-10, bad-chain-05, bad-chain-10)	275		
	E.2	Projective Dependency Parsing	276		
		E.2.1 Bilexical Unlabeled (bilexical-unlabeled)	276		
		E.2.2 Bilexical Labeled (bilexical-labeled)	277		
		E.2.3 Arc Eager (arc-eager)	277		
		E.2.4 Split Bilexical Dependency Parsing	278		
		E 2.4.1 Johnson (2007)'s Version (split-head-J)	278		

	E.2.4.2 Eisner and Blatz (2007)'s Version (split-head-EB)	278
E.3	Constituency Parsing	279
	E.3.1 CKY (cky3)	279
	E.3.2 CKY with 4-ary Productions (cky4)	279
	E.3.3 CKY with a Fixed Grammar (cky+grammar)	280
	E.3.4 Tensor Decomposition Parser (tensor-decomp-parser)	281
E.4	Inversion Transduction Grammars (itg)	281
E.5	Edit Distance (edit)	282
E.6	Bar-Hillel Construction (bar-hillel)	282
E.7	Expectations under a Linear-Chain Conditional Random Field	
	(chain-expect)	283
E.8	Intersection of Even and Odd Peano Numbers (even-odd-peano) .	284
E.9	Hidden Markov Models (hmm)	285
E.10	Semi-Markov Model (semi-markov)	285
E.11	Shortest Path in a Graph	286
	E.11.1 Explicit Path (path-list)	286
	E.11.2 Start Path (path-start)	286
E.12	Pushdown Automaton with an Explicit Stack (explicit-pda)	287
Refere	nces	290
Index		323

List of Tables

8.1	Experimental demonstration of our search-based program opti-	
	mizer on the benchmarks described in §E	224

List of Figures

3.1	Taxonomy of order relations	34
5.1	Illustration of the dimension of a derivation (Definition 29)	89
8.1	Depiction of the program optimization graph search problem. The program used in this figure is CKY (Example 6). Nodes are Dyna programs (§ 2). The node pointed to by "start" indicates the user's program. Edges are program transformations (§ 7). Costs are derived by program analysis (§ 6). Only a tiny subset of the nodes and edges in the search graph is shown. The dotted unlabeled outgoing edges represent additional transformations we did not	
	elaborate on in the diagram to reduce clutter.	208

List of Algorithms

1	unify_implicitly: Implicit unification algorithm	25
2	$\mathfrak{p}[x]$: Intersect a query x with the rules of a program \mathfrak{p} . Returns an	
	implicit result stream (i.e., an iterator that mutates x 's variable	
	bindings and emits rule bodies).	26
3	Built-in relation dispatching.	27
4	$\beta_1 \circ \beta_2$: Algorithm for joining result streams β_1 and β_2	28
5	sccs(incoming, targets): Enumerate the strongly connected compo-	
	nents (SCC; Definition 12) of a directed graph in topological order.	37
6	naive: Solve program $\mathfrak p$ by (naïve) fixpoint iteration	74
7	seminaive: Solve program $\mathfrak p$ by semi-naïve fixpoint iteration $\ \ . \ \ . \ \ .$	76
8	agenda: Solve program $\operatorname{\mathfrak{p}}$ by agenda-based fixpoint iteration	81
9	$\overline{\text{support}}$: Estimate an upper bound on the support set of $\mathfrak{p}.$	84
10	incoming and outgoing: Lazily enumerate incoming and outgoing arcs	
	in program $\mathfrak{p}\text{'s}$ dependency graph given its estimated support $\overline{\mathfrak{m}}.$.	85
11	newton: Solve program $\mathfrak p$ by Newton's method	90

12	solve_linear: Solve a right-linear program p by reduction to the
	matrix star operation (Definition 21)
13	solve: Solve program p using the hybrid algorithm
14	abstract_fc: Abstract forward chaining algorithm for type program $\overline{\mathfrak{p}}.$
	This procedure assumes that the input type has been included in $\overline{\mathfrak{p}}.108$
15	Subtype, intersection, and merge algorithms for simple types 114
16	useful: For a given program \mathfrak{p} , input type $\overline{\mathfrak{D}}$, and output type $\widehat{\mathfrak{D}}$,
	define a new set of items $\hat{\cdot}$ where \hat{x} is true if x might be useful 122
17	specialize: Specialize program $\mathfrak p$ to the type $\widehat{\mathfrak m}$
18	abbreviate: Abbreviate program $\mathfrak p$ according to the type $\widehat{\mathfrak m}$ 126
19	runtime: Infer a runtime bound for $\mathfrak p$ given a derived type $\overline{\mathfrak m}$ 136
20	unfold: Apply the basic unfold transformation (Definition 33) 162
21	fold: Apply the basic fold transformation (Definition 34) 163
22	unfold: Apply the generalized unfold transformation (Definition 35). 166
23	fold: Apply the generalized fold transformation (Definition 36) 167
24	Derivation mapping for $\mathfrak{p}^{(s)} \xrightarrow{unfold_{r \to \mathcal{S}}^{d,j}} \mathfrak{p}^{(t)}$
25	Derivation mapping for $\mathfrak{p}^{(s)} \xrightarrow{\text{fold } \mathcal{S} \to r} \mathfrak{p}^{(t)}$
26	Fold–unfold transformation sequence safety test. We assume an
	initial reference program $\mathfrak{p}^{(0)}$ with rule measure $\psi_0(\cdot)$ as input 185
27	speculation: Apply the speculation transformation (Definition 41). 191
28	not-matches: Implementation of the not-matches relation
29	leftcorner: Apply the left-corner transformation (Definition 42) 195

30	search: Algorithm for a given graph search problem $\langle \mathcal{S}, \mathcal{A}, s_0, \mathcal{T}, \text{transitions}, \text{cost} \rangle$.
	Requires a search policy π to select which candidate state to ex-
	pand. We also require that a search budget is specified, such as an
	upper limit on the total execution time spent searching 209
31	Inverse derivation mapping for $\mathfrak{p}^{(s)} \xrightarrow{unfold^{d,j}_{r \to \mathcal{S}}} \mathfrak{p}^{(t)}$
32	Inverse derivation mapping for $\mathfrak{p}^{(s)} \xrightarrow{\text{fold } \mathcal{S} \to r} \mathfrak{p}^{(t)}$
33	Derivation mapping for $\mathfrak{p}^{(s)} \xrightarrow{\operatorname{speculation}_{\boldsymbol{z},\sigma}} \mathfrak{p}^{(t)}$
34	unify: Unification algorithm. Solves any equation of the form $oldsymbol{ heta}(oldsymbol{x}) =$
	$\theta(y)$ for θ given possibly nonground terms x and y
35	cover: One-sided unification algorithm. Solves any equation of the
	form $ heta(x) = y$ for $ heta$ given possibly nonground terms x and y 267
36	covers_product: Lazily enumerate commutative matches from a col-
	lection of subgoals $\pmb{x}_1 \circ \cdots \circ \pmb{x}_N$ to another collection of subgoals $\pmb{y}_1 \circ \cdots \circ \pmb{y}_M$.
	Each iterate is a pair of a substitution mapping and an index-to-
	index alignment
37	covers_product: Lazily enumerate non-commutative matches from
	a collection of subgoals $oldsymbol{x}_1 \circ \cdots \circ oldsymbol{x}_N$ to another collection of subgoals
	$y_1 \circ \cdots \circ y_M.$ Each iterate is a pair of a substitution mapping and the
	start position of the match

Chapter 1

Introduction

Algorithmic research in natural language processing (NLP) has focused—in large part—on developing dynamic programming solutions to combinatorial problems that arise in the field. Such algorithms have been introduced over the years for numerous linguistic formalisms, such as

- finite-state transduction (e.g., Mohri (1997), Eisner (2002)),
- context-free parsing (e.g., Stolcke (1995), Goodman (1999)),
- dependency parsing (e.g., Eisner (1996), Koo and Collins (2010))
- mildly context-sensitive parsing (e.g., Vijay-Shanker and Weir (1989), Vijay-Shanker and Weir (1990), Kuhlmann, Satta, and Jonsson (2018)).

In recent years, the same algorithms have often been used for deep structured prediction, using a neural scoring function that decomposes over the structure.¹

¹E.g., Durrett and Klein (2015), Rastogi et al. (2016), Lee et al. (2016), Dozat and Manning (2017), Stern et al. (2017), Kim et al. (2017), Hong and Huang (2018), Wu, Shapiro, et al. (2018), Wu and Cotterell (2019), Qi et al. (2020), and Rush (2020).

Designing efficient algorithms is challenging. Designing and implementing dynamic programs which are as efficient as possible is challenging and error-prone. When a dynamic programming algorithm for a new problem is first introduced in the literature, its runtime is sometimes suboptimal—faster versions are often published over time (see Example 1). Indeed, introducing the first algorithm and subsequently finding improvements is common throughout computer science. We seek to facilitate this process through automation.

Example 1. Consider the following instances² of published dynamic programs whose runtime bounds were later improved. In the examples below, n denotes the length of the input sentence.

- Projective dependency parsing: Collins (1996) gave an $\mathcal{O}(n^5)$ algorithm that was sped up to $\mathcal{O}(n^4)$ by Eisner and Satta (1999a).
- Split-head-factored dependency parsing: implemented naïvely runs in O(n⁵);
 with some effort, an O(n³) algorithm can be derived (Eisner, 1996; Johnson,
 2007; Eisner and Blatz, 2007).
- Linear index-grammar parsing: $\mathcal{O}(n^7)$ in Vijay-Shanker and Weir (1989), sped up to $\mathcal{O}(n^6)$ by Vijay-Shanker and Weir (1993).
- Lexicalized tree adjoining grammar parsing: $O(n^8)$ in Vijay-Shankar and Joshi (1985), sped up to $O(n^7)$ by Eisner and Satta (2000).

²Many of these examples were brought to our attention in the works of Eisner and Blatz (2007) and Gildea (2011); further discussion can be found therein.

- Inversion transduction grammar: $\mathcal{O}(n^7)$ in Wu (1996), sped up to $\mathcal{O}(n^6)$ by Huang, Zhang, et al. (2005).
- CKY parsing (Cocke and Schwartz, 1970; Younger, 1967; Kasami, 1965)
 is typically presented in a suboptimal O(k³n³) form, but can be sped up to
 O(k²n³ + k³n²) (Lange and Leiβ, 2009; Eisner and Blatz, 2007) where k is the
 number of nonterminal symbols in the grammar.
- Tomita's context-free parsing algorithm (1985) runs in O(n^{ρ+1}) where ρ is the length of the longest right-hand side of a context-free production in the grammar (Johnson, 1989). However, it can be made to run in O(n³) by binarizing the production rules.

Central question. This dissertation explores the question:

Can we automatically discover these faster algorithms?

A unified notation for dynamic programming. Our approach starts with an observation made by Goodman (1999) that a dynamic programming algorithm can be regarded as performing inference in a semiring-weighted deduction system. Their work builds upon a rich history of (unweighted) deduction systems that we review in §9.1. Eisner, Goldlust, et al. (2005) provided a programming language, Dyna, for expressing such deduction systems, along with a compiler that produced fast inference code. We will use Dyna in this thesis as our algorithm-specification language. We introduce the language in §2 and give formal semantics in §4.

Below is an example Dyna program that computes the total weight of all paths between pairs of nodes in a weighted directed graph.

```
path(I,I) += 1.
path(I,K) += edge(I,J) * path(J,K).
```

Here the weight of an edge comes from the relation edge(I,J). For example, to declare that the weight of the edge from a node labeled "a" to a node labeled "b" is 13, we write edge("a", "b") += 13. From there, the path(I,K) items encode a weighted relation between all the nodes in the graph. Specifically, for all pairs of nodes I,K:

$$path(I,K) = \underbrace{1[I=K]}_{\#_I} + \sum_{J} \underbrace{edge(I,J) \cdot path(J,K)}_{\#_2}$$
(1.1)

Here the variable J ranges over the set of nodes. The reason why the variable J is summed in (1.1), while I and K are not, is because it is local to the right-hand side of #2. The reader may recognize (1.1) as the mathematical recurrence for the total weight of all paths between pairs of nodes in a graph. Turning this mathematical specification into an efficient computational procedure requires several implementation details to be specified. In this example, we need to select data structures to store and index the edge and path relations—but, more importantly, we need to *solve* the recursive equation to derive the total weight. In §5, we describe general-purpose procedures for taking the mathematical recurrences (written in Dyna) and turning them into efficient algorithms.

Automated analysis. A key ingredient to our system is static analysis: the ability to reason about the program without running it on specific inputs.

Such analyses include type analysis, runtime complexity analysis, and space complexity analysis (Cousot and Cousot, 1977; McAllester, 2002). These analyses are important for understanding the space of possible program behaviors on unknown input data. Automated program analysis aids programmers because manual program analysis may be tedious and error-prone. In the context of NLP, automated program analysis of dynamic programs (for instance, as specified in Dyna) has the potential to help NLP algorithms researchers correctly and speedily analyze their creations. Indeed, the NLP literature provides examples where such automated analysis could have been of use.

Example 2. Consider the following example from the NLP literature:

- Kuhlmann, Gómez-Rodríguez, et al. (2011, §5.4) discovered a surprising property of their tabular dependency parsing algorithm. Their initial naive analysis gave a time complexity of $\mathcal{O}(n^5)$ and space complexity of $\mathcal{O}(n^3)$. Still, they were able to tighten this analysis to a time complexity of $\mathcal{O}(n^3)$ and space complexity of $\mathcal{O}(n^2)$ after noticing that the objects built by their transition system always had two equal variables (which could be proved inductively). They then revised their program formulation to eliminate the redundant variable.
- Huang and Sagae (2010) derive a projective dependency parsing algorithm with a stated complexity of $\mathcal{O}(n^7)$, which was later discovered to actually be $\mathcal{O}(n^6)$ according to Shi et al. (2017).⁴

³Although we do not explore it in this thesis, we also mention that static program analysis is essential for compiling efficient code. Type inference, for example, allows the compiler to synthesize efficient data structures for representing internal objects and efficient code for pattern-matching against them.

⁴This is stated in Shi et al. (2017) without proof. The only explanation we could find in the

Our automated analysis system yields big- \mathcal{O} bounds on runtime and space complexity (§6.3). These bounds are derived by inferring types for each variable in a given rule and using each type to bound the number of ways the rule will fire during execution (§6.3). At a high level, to apply this technique to the path program, we first see that the variables I, J, and K are nodes. Supposing that the set of nodes has size n. The number of ways that the rule $\operatorname{path}(I,K) := \operatorname{edge}(I,J) : \operatorname{path}(J,K)$ will fire during inference is $\mathcal{O}(n^3)$. The foundations of this technique are given by McAllester (2002)'s meta theorem. Many authors⁵ manually apply the meta theorem to analyze the runtime of their recursive algorithms. Our runtime analyzer (§6.3) automatically applies the meta theorem to derive tight bounds on the worst-case space and runtime complexity of executing a given program (using the execution methods in §5).

Systematizing algorithmic speed-ups. All of the runtime improvements mentioned in Example 1 can be cast as source-to-source program transformations on Dyna programs (Eisner and Blatz, 2007); we develop these transformations in §7.6 In dynamic programming, program transformations may be exploited to

literature is in Hong and Huang (2018). In that work, the authors explain that one of the variables in the reduction rule (of a closely related model), which appears to take on $\mathcal{O}(n)$ values, is functionally dependent on the other variables. This means that it may take on at most one value given the values of other variables in scope; thus, one of the factors of n in the runtime becomes a constant. This observation appears to also hold in the model of Huang and Sagae (2010).

⁵A list of NLP papers that explicitly use this analysis technique includes the following citations: Gildea (2011), Nederhof and Satta (2011), Gilroy et al. (2017), Melamed (2003), Kuhlmann (2013), Nederhof and Sánchez-Sáez (2011), Büchse et al. (2011), Lopez (2009), Eisner and Blatz (2007), Huang and Sagae (2010), and Kuhlmann, Gómez-Rodríguez, et al. (2011).

⁶In a seminal paper, Burstall and Darlington (1977) introduced the idea of synthesizing faster/better (functional) programs by source-to-source transformation: given an initial simple, lucid, and—most importantly—correct program, we may repeatedly transform it by altering its recursive structure to eventually reach an improved version of the program. Later, Tamaki and Sato (1984) extended the technique to logic programming.

derive algorithms with a faster runtime. These transformations map a program to another with the same meaning (i.e., given the same inputs, it will produce the same outputs) but with a possibly different running time. An example of such a speed-up is given below.

Suppose that we only want to compute the total weight of paths from one node set start(I) to another node set stop(K). The following program does this by first computing the total weights for all pairs and then summing over a subset:

```
goal += start(I) * path(I,K) * stop(K).
path(I,I) += 1.
path(I,K) += edge(I,J) * path(J,K).
```

However, this program is not as efficient as it can be. By exploiting the distributive property of addition and multiplication, we can sum over the possible stop nodes sooner in the recurrence:

```
goal += start(I) * pathstop(I).
pathstop(I) += stop(I).
pathstop(I) += edge(I,J) * pathstop(J).
```

Here pathstop(I) computes the total weight of all paths from node I to a node in stop. The running time of this program is $\mathcal{O}(n^2)$. In §7, we provide a collection of program-to-program transformations capable of recovering this speedup.

Program improvement through automated transformation search.

This work shows how to search over program transformation sequences to automatically discover faster algorithms, automating the work of the NLP algorithmist. Our work poses program optimization as search over transformed versions of the initial program. We describe how to use search to rapidly rediscover many of the known optimizations listed in Example 1. The program optimizer ties

together many aspects of the thesis: Our space of possible programs (§2), a cost function that serves as a proxy for program runtime (§6), and a set of directed edges that connect semantically equivalent programs (§7). Our search starts at the initial program and seeks a low-cost equivalent program that can be reached by traversing directed edges. In §8, we review the search techniques that we will use for this purpose. We provide an experimental demonstration in §8.3.

Summary. We motivated our approach with many examples from NLP illustrating that designing and analyzing efficient dynamic programs is challenging. We will develop a framework that decomposes the problem of systematically analyzing and improving dynamic programs into the following components:

- 1. a domain-specific language to logically specify an algorithm (§ 2 and § 4) and general machinery to execute them (§ 3.1 and § 5)
- 2. a system for the automated analysis of the space/runtime complexity of dynamic programs as well as type analysis (§6)
- 3. a rich collection of meaning-preserving transformations to programs (§7)
- 4. a search algorithm for identifying a good sequence of transformations that reduce the running time complexity (§8).

We deliver a powerful system that can recover the kinds of speedups in Example 1.⁷ We present an experimental demonstration of our system in §8.3. A prototype implementation is freely available at https://github.com/timvieira/dyna-pi.

⁷Our experiments include all of the speedup scenarios of Example 1 *except* for the linear index-grammar and lexicalized tree adjoining grammar parsers. We include many other interesting speedup scenarios in our experiments that we have not mentioned here.

Chapter 2

A Domain-Specific Language for

Dynamic Programming

Dyna is a high-level, domain-specific language for concisely expressing dynamic programs as recurrence relations. It was first introduced by Eisner, Goldlust, et al. (2005), where it was demonstrated to be a useful notation and software library for rapidly developing efficient algorithms for a variety of popular models of linguistic structure (such as tagging, parsing, and translation) as well as many combinatorial optimization problems that can be solved by dynamic programming (such as shortest path, knapsack, and sequence alignment).

Our exposition of the Dyna language will unfold as a sequence of simple and hopefully familiar computations—introducing concepts and terminology when necessary or convenient.

Example 3. The total weight of length-4 paths in a weighted directed graph:

```
9 total += w(Y_1, Y_2) * w(Y_2, Y_3) * w(Y_3, Y_4) * w(Y_4, Y_5).
```

This program defines the **value** of a derived **item** total in terms of input items w(...). This program assumes that we have encoded the set of edges in the given graph by declaring values for w(...) items. For example, to declare that there is an edge from a node labeled "a" to a node labeled "b" with weight 13, we write w("a","b") += 13. The value of total is $\sum_{Y_1} \cdots \sum_{Y_5} w(Y_1,Y_2) \cdot w(Y_2,Y_3) \cdot w(Y_3,Y_4) \cdot w(Y_4,Y_5)$.

Dyna rules may be recursive:

Example 4. Minimum-cost path in a graph.

```
goal min= start(I) + \beta(I).

\beta(J) min= stop(J).

\beta(I) min= w(I,J) + \beta(J).
```

This program defines $\beta(I)$ as the cost of the minimum-cost path in a graph from the node I to any terminal node J. Here stop(J) must first be defined to be the cost of stopping at J when J is a terminal node (and should be left undefined otherwise). Also, w(I,J) must first be defined to be the cost of the directed edge from I to J (and should be left undefined if there is no such edge).

Separation of logic and control. Example 4 specifies the problem of finding the single-target shortest path in a weighted directed graph. However, the details of the implementation are abstracted away. Thus, many well-known single-target shortest path algorithms, such as the Bellman–Ford–Moore algorithm (Bellman, 1958; Ford, 1956; Moore, 1959) and Dijkstra's algorithm (Dijkstra, 1959) are represented as the same Dyna program (Example 4) as they are based on the same recursive factoring of paths. The differences between

these algorithms are lower-level than the basic Dyna notation allows as they are essentially only differences in evaluation order (§5.1.3). One of the benefits of the Dyna notation is that it separates the mathematical specification of the algorithm from the details of the execution of that algorithm (e.g., data structures, control flow, execution order). This idea is often referred to as *the separation of logic and control* (Kowalski, 1979).

Notice that if we change the equations from $\langle \min, + \rangle$ to $\langle \max, * \rangle$, we get the Viterbi algorithm (Viterbi, 1967) for finding the most probable path in a graph from nodes in a set start to nodes in a set stop where the probability of a path is given by the product of edge probabilities (or weights) w(...):

Example 5. Finding the most probable path in a graph.

```
goal max= start(J) * \beta(J).

\beta(J) max= stop(J).

\beta(J) max= \beta(I) * \beta(I, J).
```

What is interesting about the connection between Example 4 and Example 5 is that we kept the same recurrence, but we have changed the aggregate quantity that we seek to compute. Both programs define the same set of derivations, which we can interpret as paths in the graph. How these paths are weighted and interpreted differs. In Example 4, paths have costs; the costs are extended with addition and aggregated with min. In Example 5, paths have probabilities; the probabilities are extended with multiplication and aggregated with addition. We will return to the extension and aggregation operations momentarily.

Next, we give a fundamental algorithm to NLP.

Example 6. Weighted context-free parsing with CKY (Cocke and Schwartz, 1970; Younger, 1967; Kasami, 1965), or, more precisely, the inside algorithm (Baker, 1979; Jelinek, 1985):⁸

```
16 \beta(X,I,K) += \gamma(X,Y,Z) * \beta(Y,I,J) * \beta(Z,J,K).

17 \beta(X,I,K) += \gamma(X,Y) * \beta(Y,I,K).

18 \beta(X,I,K) += \gamma(X,Y) * word(Y,I,K).

19 goal += \beta(S,0,N) * len(N).
```

The values of the γ items should be defined to be the weights of the corresponding context-free grammar rules: for example, the item $\gamma(s,np,vp)=0.7$ encodes the production $s \xrightarrow{0.7} np vp$. Also, the item word(X,I,K) should be 1 if the input word χ appears at position κ of the input sentence and χ and should be left undefined otherwise. Then for any nonterminal symbol χ and any substring spanning positions (I,K] of the input sentence, the item $\beta(\chi,I,K)$ represents the total weight of all grammatical derivations of that substring from χ .

Syntax and Terminology. A Dyna **program** \mathfrak{p} serves to define **values** for **items**. Thus, the program defines a function from items to values; we denote the value of an item x as [x]. We call $[\cdot]$ the **valuation function**; we defer its formal definition to Definition 22. Each item is named by an element of the **Herbrand universe** \mathbb{H} —that is, a (possibly nested) labeled tuple such as f(g(z,h(3))), which is a 1-tuple labeled f whose single element is a 2-tuple labeled g. The elements of f are known as **ground terms** because they contain no variables. The program f itself is an unordered collection of **rules** of the form f of the form f of the body is called a f the **head**, and f and f of the rule. Each f in the body is called a

⁸If the reader is not familiar with context-free parsing, we recommend Jurafsky and Martin (2020, chapters 12–13).

⁹Note that the collection may contain duplicates (§3.1.1).

subgoal.

The operations \oplus and \circ manipulate 10 values from a set $\mathbb V$ and may be any pair of operations that form a **semiring** (§3.3). Let $\mathbb Q$ and $\mathbb I$ denote the semiring's additive and multiplicative identity elements, respectively. If an item's value is left undefined or has no \oplus -summands, this is equivalent to saying it has value $\mathbb Q$, which means it cannot affect the values of any other items; in effect, it does not exist. Common semirings include total weight $(\mathbb R_{\geq 0}, +, *)$, boolean $(\{\text{false}, \text{true}\}, \vee, \wedge)$, minimum $\text{cost} (\mathbb R \cup \{\infty\}, \text{min}, +)$, and Viterbi $(\mathbb R_{\geq 0}, \text{max}, *)$. The distributive property of the semiring gives rise to an important interpretation: The value of a given item is a certain sum over all of its **derivations** (i.e., proof trees) under the program's rules. The value of the derivation is the product of the values of its axioms. In this paper, we assume that all rules in the program use the same semiring.

Rules do not have to name specific items. The terms in a rule may include **variables**—denoted by capitalized letters, such as x—that are universally quantified over \mathbb{H} . This allows them to pattern-match against many items. Let x_1, \ldots, x_N denote the distinct variables in a rule R. A **grounding** of R is a

¹⁰The notation used to denote the semiring operations varies from author to author. Many authors (Lehmann, 1977; Kuich, 1997; Green, Karvounarakis, et al., 2007; Esparza et al., 2007) use $\langle +, \cdot, 0, 1 \rangle$, which leads to confusion since these operations are often used to denote their ordinary meanings in the same text. In the NLP literature, it is most common to $\langle \oplus, \otimes, \odot, \odot \rangle$ to avoid excessive overloading of $\langle +, \cdot, 0, 1 \rangle$. We have chosen to use $\langle \oplus, \circ, \underline{0}, \underline{1} \rangle$ since they are easier to distinguish visually.

¹¹Semirings provide an elegant abstraction for computing many common quantities in NLP models, such as expected values, covariance, gradients, entropies, and *K*-best (Goodman, 1999; Huang, 2008; Eisner, Goldlust, et al., 2005; Huang and Chiang, 2005; Li and Eisner, 2009).

¹²The precise semantics, including the definition of the set of derivations, is given in §4.

¹³We formalize the sum-over-derivations interpretation in §4.

 $^{^{14}}$ Dyna $\underline{2}$ (Eisner and Filardo, 2011), relaxes this restriction. We discuss future work on Dyna $\underline{2}$ in $\S D$.

variable-free rule $\{X_1 \mapsto X_1, \dots, X_N \mapsto X_N\}(\mathbf{R})$ (that is, the result of substituting x_i for all copies of X_i in \mathbf{R} , for each i) for some $x_1, \dots, x_N \in \mathbb{H}$. $\Gamma(\mathbf{R})$ denotes the **set of all groundings** of \mathbf{R} ; including a non-ground rule in \mathfrak{p} is equivalent to including all of its groundings.

The Dyna language allows logical **side conditions** on a rule, such as goal += f(X) for X < 9. This is syntactic sugar for goal += f(X) * lesstham(X,9), where the value of the lesstham constraint is given by

$$[lessthan(X,Y)] \stackrel{\text{def}}{=} \underline{1} \text{ if } X < Y \text{ else } \underline{0}$$

The constraint x < 9 makes use of the **built-in relation** <; Dyna's standard library supports many other common arithmetic relations, such as

- $[X \text{ is } Y+Z] \stackrel{\text{def}}{=} \underline{1} \text{ if } X=Y+Z \text{ else } \underline{0}$
- $[X \text{ is } Y^K] \stackrel{\text{def}}{=} \underline{1} \text{ if } X = Y^K \text{ else } \underline{0}$

Unlike the items defined by the program's rules, built-in relations are opaque—they are not defined by rules but rather by executable (non-Dyna) code.

Input data is represented as program \mathfrak{D} consisting of **axioms**—i.e., rules of the form $(x \oplus v)$ where $v \in V$. For example, consider the shortest-path program as \mathfrak{p} and the specific graph to run it on as \mathfrak{D} .

$$\mathfrak{p} = \begin{cases} 20 & \text{\% definition of shortest} \\ 21 & \text{\% path from S to stop} \\ 22 & \beta(S) \text{ min= w(S,S')} + \beta(S'). \\ 23 & \beta(S) \text{ min= stop(S)}. \end{cases} \mathfrak{D} = \begin{cases} 24 & \text{\% define graph (w, stop)} \\ 25 & \text{w(bal,phl) min= 3.} \\ 26 & \text{w(phl,nyc) min= 1.} \\ 27 & \dots \\ 28 & \text{stop(nyc) min= 0.} \end{cases}$$

Adding the $\mathfrak D$ rules to $\mathfrak p$ results in a new program $\mathfrak p'$ that can subsequently be

queried, as we will show below.

Querying a Program. The way that a user interacts with a Dyna program is through queries. A **query** against program \mathfrak{p} is specified as a (possibly nonground) term x, and the **result** of a query is a **specialized program** that encodes the restriction of the valuation function $[\cdot]$ (Definition 22) to the possible groundings of x. Below is an example of a program, a query, and the query's results (going left to right).

```
\begin{cases} 29 & \text{\% definition of shortest} \\ 30 & \text{\% path from S to stop} \\ 31 & \beta(S) & \text{min=} \ w(S,S') + \beta(S'). \\ 32 & \beta(S) & \text{min=} \ \text{stop}(S). \\ 33 & \text{\% define graph (w, stop)} \\ 34 & \text{w(bal,phl) min=} \ 3. \\ 35 & \text{w(phl,nyc) min=} \ 1. \\ 36 & \dots \\ 37 & \text{stop(nyc) min=} \ 0. \end{cases} \qquad \underbrace{\frac{\text{query }\beta(S)}{\text{query }\beta(S)}}_{\text{query }\beta(S)} \qquad \underbrace{\frac{38}{39} \ \beta(\text{bal) min=} \ 4.}_{\text{39} \ \beta(\text{phl}) \ \text{min=}} \ 1. \\ 40 & \beta(\text{nyc}) \ \text{min=}} \ 0. \end{cases}
```

The query $\beta(S)$ returns a map: each item of the form $\beta(S)$ maps S to the cost of its shortest path to any of the Stop nodes. As the example output above shows, Dyna represents the result of a query as specialized programs, i.e., programs that contain only the answers to the specific query. The specialized programs that are answers to queries are restricted to have the property that all their right-hand sides are constants. ¹⁵

Notice that the results for a query may be an infinite-sized valuation. To compactly represent certain infinite valuations, we allow the heads of the rules in the results to contain free variables. For example,

¹⁵It is possible to relax this condition to allow the right-hand side of a rule to contain built-in relations. For example, f(X) += 3 for (1 is X mod 2) says that f(X) items have a contribution of 3 if X is an odd number.

For efficiency, we apply **constant folding** to simplify the contribution to each head wherever possible. For example, the following results would be replaced by #42-43.

```
44 f(X) += 6.

45 f(Y) += 1.

46 f(3) += 9.

47 f(3) += 2.
```

Summary

This chapter has provided a high-level overview of the Dyna syntax. We introduced the main concepts and terminology Dyna users should be familiar with. We will describe the concepts pertaining to the internal representation of programs in §3, the precise meanings of a derivation and the valuation mapping in §4, and how to execute programs in §5.

Chapter 3

Notation, Background, and

Utilities

This chapter includes background material that we will use throughout this thesis. This background material includes:

3.1	Interr	nal Representation of Programs	18
	3.1.1	Notation and Representation	18
	3.1.2	Lookup and Update Represention	22
3.2	General Mathematics		31
	3.2.1	Ordering Relations	32
	3.2.2	Functions	34
	3.2.3	Graphs	35
	3.2.4	Abstract Algebra	38
3 3	Semirings: A Powerful Abstraction for Dynamic Programming		49

3.1 Internal Representation of Programs

This section describes important notation and our internal program representation. We will describe an efficient data structure for representing the results of a query as a *result stream* (§3.1.2.1). We will develop a representational abstraction that allows us to regard rules as variable-free when manipulating them (e.g., during inference (§5)).

3.1.1 Notation and Representation

We will now describe how programs are represented internally. When the user writes a program such as

```
48 f(X) += g(X,Y) * h(Y).
49 f(X) += 1.
```

it is represented as a bag of rules. Each rule R is of the form $x \oplus = y_1 \circ \cdots \circ y_K$. It is represented internally as a pair of a head and a body. The body is a list of subgoals. Let head(R) and body(R) denote the head and body of the rule R. The head x and each subgoal y_k are represented as nonground terms. A **nonground term**¹⁶ is a labeled tuple that may contain named variables (denoted by capital letters, such as x and y in the example). Let $vars(\cdot)$ be a function that returns the set of free variables contained in a nonground term, e.g., $vars(f(g(x), 4, x)) \mapsto \{x\}$. Variables can be shared between nonground terms, e.g., x is shared between the head f(x) and the subgoal g(x, y) of #48 and y is shared between its two subgoals.

¹⁶When we write "nonground term," we mean a *possibly* nonground term, i.e., a term that is allowed to be nonground but does not necessarily contain free variables.

Substitution is a function $\theta(\cdot)$ from nonground terms to nonground terms; it is parameterized by a **substitution mapping** θ that maps variables to nonground terms, such as $\theta = \{Y \mapsto g(X), Z \mapsto 4\}$. When we apply the substitution map, $\theta(x)$, it will replace all instances of the variable on the left-hand side of an arrow (\mapsto) with whatever is on the right-hand side of the arrow.

Unification. The primary operation performed on nonground terms is **unification**; it is performed internally by our solvers (§ 5) to bind variables. Unification solves equations of the form $\theta(x) = \theta(x')$ where x and x' are nonground terms. Such constraints can be solved efficiently using a unification algorithm, $\theta \leftarrow \text{unify}(x,x')$ (e.g., Robinson (1965) and Martelli and Montanari (1982)).¹⁷ Here unify returns a substitution mapping θ or FAIL if no solution exists. For example, $\theta \leftarrow \text{unify}(f(Y,4), f(g(X),Z))$ returns the solution $\theta = \{Y \mapsto g(X), Z \mapsto 4\}$. When we apply this substitution map to both sides, the terms become equal, $\theta(f(Y,4)) = \theta(f(g(X),Z))$. On the other hand, $\theta(f(X,g(X))) = \theta(f(3,g(4)))$ is unsatisfiable; thus, $\theta = \text{FAIL}$. We also mention equations of the form $\theta(x) = x'$. Such constraints occur in program transformations §7, and they can be solved in a similar fashion using $\theta \leftarrow \text{cover}(x,x')$; pseudocode is given in Algorithm 35.

Since variable names are global, we often need to rename them (the x in #48 is supposed to be different from the x in #49. For this purpose, we use a function fresh, which renames all variables within a term, rule, or entire program to novel names in order to avoid unwanted name conflicts. For example,

¹⁷We give pseudocode and implementation details for our unification algorithm in §C.1.

¹⁸We allow unify to take a substitution mapping as an optional third argument, which enables algorithms to chain together the results of unification calls.

 $fresh(f(g(X),X)) = f(g(X_*),X_*)$ where X_* is a novel variable name.

Destructuring. We will use destructuring, such as $(x \oplus = y_1 \circ \cdots \circ y_K) \leftarrow \mathbf{R}$, to break up the rule \mathbf{R} into its head and variable-length body. It is equivalent to $x \leftarrow \text{head}(\mathbf{R})$, $K \leftarrow |\text{body}(\mathbf{R})|$, $\langle y_1, \dots, y_K \rangle \leftarrow \text{body}(\mathbf{R})$. If K = 0, we regard the body as $\underline{\mathbf{1}}$. We may write $f(x_1, \dots, x_N) \leftarrow x$ to destructure a term into its arguments and bind f to its label. For example, $f(x_1, \dots, x_N) \leftarrow \mathbf{p}(\mathbf{q}(X), \mathbf{r}(X))$ sets $N \leftarrow 2$, $f \leftarrow \mathbf{p}$, $x_1 \leftarrow \mathbf{q}(X)$, $x_2 \leftarrow \mathbf{r}(X)$. In a destructuring expression, such as $(\underline{} \oplus y_1 \circ \cdots \circ y_K) \leftarrow \mathbf{R}$, we may use underscore $(\underline{})$ to match a subexpression that we do not wish to name.

Bags. We will use bags (aka multisets) to represent a collection of rules (programs). Literal bags are denoted by the paired delimiters $\langle \cdot \rangle$. Let α and β be bags. We write operations $\alpha \cup \beta$ for bag union and γ for bag difference. When used as a condition, $x \in \beta$ is true if and only if x appears at least once in β . However, when used in a loop, $x \in \beta$ will repeat duplicate instances of x. We will access elements of β_i of β by their "position" i—even though the position is arbitrary—these positions are simply a canonical linear order of the (possibly duplicated) elements.

Programs. The notation $\mathfrak{p}[x] \oplus \mathfrak{p}_1 \circ \cdots \circ \mathfrak{p}_K$ adds the rule $x \oplus \mathfrak{p}_1 \circ \cdots \circ \mathfrak{p}_K$ to the program \mathfrak{p} . We use the notation $\mathfrak{p}[x]$ to retrieve the collection containing the bodies of the (specialized) rules in \mathfrak{p} whose heads have been unified with the subgoal x. To take the product of two subgoals x and x', we may write $\mathfrak{p}[x \circ x']$ or $\mathfrak{p}[x] \circ \mathfrak{p}[x']$. The notation can be used to read (rule bodies) and write (new rules)

 $^{^{19}\}mathrm{We}$ may optionally pair the rule body with the head-unificiation's substitution mapping.

across programs with expressions like $\mathfrak{p}'[x] \oplus \mathfrak{p}[y_1] \circ \cdots \circ \mathfrak{p}[y_K]$. We write $\mathfrak{p} \oplus \mathfrak{p}'$ to concatenate the rules of two programs: \mathfrak{p} and \mathfrak{p}' . To create an empty program \mathfrak{p} , we write $\mathfrak{p} \leftarrow \emptyset$. We explain all of this notation in more detail in §3.1.2.

Shortcuts for defining integer intervals. We find the following notation to be useful in defining lists of contiguous integers, $i, j \in \mathbb{Z}$:

$$\begin{split} &[i:j) \stackrel{\text{def}}{=} \left[\ i,i+1,\ldots,j-1 \quad \right] \\ &[i:j] \stackrel{\text{def}}{=} \left[\ i,i+1,\ldots,j-1,j \ \right] \\ &(i:j) \stackrel{\text{def}}{=} \left[\quad i+1,\ldots,j-1 \quad \right] \\ &(i:j) \stackrel{\text{def}}{=} \left[\quad i+1,\ldots,j-1,j \ \right] \end{split}$$

When an endpoint is not specified, we take it to be an extreme value, e.g., [-3:) = [-3, -2, ...] and [:10) = [..., 9, 10]. We use this notation primarily to index rule bodies.

Rule bodies. We make use of the following notational shortcuts for manipulating rule bodies. Let $\mu = y_1 \circ \cdots \circ y_K$ be a list of subgoals. Let β be a list of indices, then we define $\mu_{\beta} \stackrel{\text{def}}{=} [y_i \mid i \in \beta \cap [1:K]]$. We frequently use the integer-range syntax for β , e.g., $\mu_{[i:j)}$ denotes the contiguous range of subgoals in [i:j). We may also write $y_{[i:j]} \circ y_{(j:k]} = y_{[i:k]}$. Similarly, $\mu_{[i:j)} \circ \mu_{j} \circ \mu_{(j:k]} = \mu_{[i:k]}$, and $\mu_{[i:j)} \circ \mu_{j} \circ \mu_{(j:l)} = \mu$.

²⁰In these extreme cases, the difference between an open ("(" or ")") and closed ("[" or "]") endpoint is irrelevant.

3.1.2 Lookup and Update Represention

In this section, we will develop a notation that substantially reduces the clutter that results from explicitly passing around substitution mappings and from iteration over the results of queries. On the left is the notation that we will develop, and on the right is the code that it represents.

$$\mathfrak{p}[\mathsf{a}(\mathtt{I},\mathtt{L})] \oplus = \\ \mathfrak{m}[\mathsf{b}(\mathtt{I},\mathtt{J})] \circ \mathfrak{m}[\mathsf{c}(\mathtt{J},\mathtt{K})] \circ \mathfrak{m}[\mathsf{d}(\mathtt{K},\mathtt{L})] \\ \mathfrak{m}[\mathsf{b}(\mathtt{I},\mathtt{L})] \circ \mathfrak{m}[\mathsf{c}(\mathtt{J},\mathtt{K})] \circ \mathfrak{m}[\mathsf{d}(\mathtt{K},\mathtt{L})] \\ \mathfrak{m}[\mathsf{b}(\mathtt{I},\mathtt{L})] \circ \mathfrak{m}[\mathsf{c}(\mathtt{J},\mathtt{L})] \circ \mathfrak{m}[\mathsf{d}(\mathtt{K},\mathtt{L})] \\ \mathfrak{m}[\mathsf{b}(\mathtt{I},\mathtt{L})] \circ \mathfrak{m}[\mathsf{c}(\mathtt{J},\mathtt{L})] \circ \mathfrak{m}[\mathsf{d}(\mathtt{K},\mathtt{L})]$$

In the example, there are two programs: \mathfrak{p} and \mathfrak{m} . We are looking up values in \mathfrak{m} , and we are updating \mathfrak{p} . The bracket notation $\mathfrak{m}[x]$ returns a stream of pairs of values and bindings to the variables in x.²¹ The notation $\mathfrak{p}[x] \oplus = \beta$ will add rules to \mathfrak{p} that result from expanding the queries in β .²² We will formalize the lookup and update notation shortly.

An **explicit result stream** for that query x represents the results of a query against a given program as a stream: $[\theta^{(t)}((x \oplus v^{(t)}))]_{t=1...}$ where each $\theta^{(t)}$ is a substitution map (over x's variables) and $v^{(t)} \in V$ is a contribution to the value of the item $\theta^{(t)}(x)$. We will employ one more trick: we will *implicitly* update the bindings of variables in x as our index t in the stream progresses. We do so by representing variables as pointers and the substitution mapping by what each

²¹The notation $\mathfrak{m}[x]$ should not be confused with $\mathfrak{m}[x]$, which computes the valuation mapping for items matching x.

²²Note that the result of expanding these queries could be an arbitrary rule body.

²³The result streams have the added benefit that laziness supports infinitely long results. Such flexibility can be useful if the stream's consumer will not consume it entirely (i.e., the infinite stream is *short-circuited*).

variable (pointer) points to. We call this an **implicit result stream**. This has the benefit that if a variable in x occurs outside of x, its bindings are seamlessly and efficiently propagated. For example, if we create a result stream for f(X), then the bindings to X would be propagated to g(X) because they share the variable X. This implementation strategy has the benefit that $\theta^{(t)}$ —and substitutions by it—can be abstracted away in most algorithms—allowing for concise pseudocode that does not need to directly manipulate or manage $\theta^{(t)}$.

3.1.2.1 Implicit Result Stream Data Structure

Our implicit result system implements substitution by directly binding the variables in x through in-place unification (Algorithm 1). Each variable is represented as a pointer (or reference) in this system. Whenever distinct variables x and y are unified, we can simply bind the x pointer to y. Later, if we want to test the equality of two variables, we follow each of their pointers to their respective roots and compare literal equality. We extend $vars(\cdot)$ to return the set of variables under the control of a given result stream.

Since our iterators are in-place, in order to proceed to iteration t+1, the iterator needs to *undo* the unifications from iteration t before it can safely apply any new unifications. The in-place implementation avoids unnecessary copying operations (e.g., calls to fresh, copying θ). To facilitate undoing a unification, we

²⁴A similar data structure underlies Warren's Abstract Machine, which is the basis of most Prolog implementations (Warren, 1983; Aït-Kaci, 1991).

²⁵Our data structure is inspired by the well-known union–find data structure (Galler and Fischer, 1964) for representing equivalence relations. The use of union–find for unification was explored by Mannila and Ukkonen (1986).

have implemented unify_implicitly as a Python-style iterator, which has localized *do* and *undo* code. If unification fails, this iterator loops zero times. If unification succeeds, the iterator modifies the variable bindings as a side effect and yields a value, then suspends while this result is processed. Once this iterator resumes, it will undo its bindings—restoring its arguments to their prior state from before unify_implicitly was called.

Pseudocode for implicit unification is provided in Algorithm 1; it uses syntax for creating iterators based on Python's generators. The logic in the procedure is the same as unify (Algorithm 34). The pseudocode makes use of the following helper functions: isbound(\cdot), which returns true if and only if its argument is not a null pointer (i.e., not a free variable), and $deref(\cdot)$ follows the pointer.

Lastly, we write [x = y] as a shorthand for unify_implicitly (x, y), 27 and v[x = y] as shorthand for $v \circ [x = y]$ where $v \in V$ is constant value.

3.1.2.2 Lookup (p[x])

In this section, we show how to create (implicit) result streams by intersecting a query with the rules of a program and joining result streams so that they correspond to product queries. Pseudocode for the intersection of a query x against a program p is given in Algorithm 2. In this procedure, the query x is

²⁶Python's generator syntax (https://peps.python.org/pep-0255/) includes the **yield** keyword, similar to **return** in regular functions. However, the key difference is that **return** provides a *single* value, whereas **yield** allows for the production of a *sequence* of values. Moreover, generators can be paused and resumed: a **yield** statement outputs a value and temporarily suspends execution. When the generator is next called, it resumes from where it left off, continuing until it encounters another **yield** or it reaches the end of the function. Additionally, the expression "**yield from** xs" is a shorthand for "**for** $x \in xs$: **yield** x".

²⁷This notation was inspired by Iverson bracket notation.

Algorithm 1 unify implicitly: Implicit unification algorithm.

```
1. def unify implicitly (x, y):
     if x = y: yield 1
                                                                                      \triangleright Success: suspend
      else if isvar(x):
         if isbound(x): yield from unify implicitly(deref(x), y)
         else if isvar(y) and isbound(y): yield from unify_implicitly(x, deref(y))
         else if x does not occur in y:
            try:
               bind x to y
                                                                              \triangleright x is free; bind it to y here
               yield 1
                                                                                      ▷ Success: suspend
            finally:
               unbind x

▷ Undo binding

11
      else if isvar(y): yield from unify_implicitly(y, x)
                                                                                \triangleright Handled by symmetry
      else if isterm(x) and isterm(y) and their labels are equal:
13
         yield from unify_implicitly(x.args, y.args)
14.
      else if isseq(x) and isseq(y) and |x| = |y|:
                                                                           □ Unify argument sequences
15
         if |x| = 0: vield 1
                                                                                      ▷ Success: suspend
16
         else:
17
            {f for}\ \_\in {\sf unify\_implicitly}({\it x}_{\!1},\,{\it y}_{\!1}) :
18
               \mathbf{yield} \; \mathbf{from} \; \mathsf{unify\_implicitly} \Big( \pmb{x}_{(1:]}, \pmb{y}_{(1:]} \Big)
19
```

nondeterministically matched against different cases, including values, builtins, and program rules. In the case of program rules, if the query x unifies with the head x' of a rule $x' \oplus y_1 \circ \cdots \circ y_K$, the dispatch procedure suspends—allowing the caller of the method to peer at the state of the variables in x and the value $y_1 \circ \cdots \circ y_K$ emitted by **yield**. Once the caller is ready to move to the next match, the process continues to the following rules in p. Notice that this procedure does not compute p[x]; it only expands the query x by a single step.

Built-in relations. On line 2.6, we indicate that built-in relations should insert their logic into the lookup procedure. To give a sense of what that logic

²⁸In practice, we speed-up the loop line 2.8 by indexing the rules of $\mathfrak p$ by a "hash code" that ensures that for any q, $\mathsf{unify}(\mathsf{head}(s), q) \neq \mathsf{FAIL} \Rightarrow \mathsf{hash}(\mathsf{head}(s)) = \mathsf{hash}(q)$. Additionally, the cost of fresh can be amortized somewhat by keeping a pool of fresh copies for each rule and reusing a copy from the pool if there are no active references to it.

Algorithm 2 $\mathfrak{p}[x]$: Intersect a query x with the rules of a program \mathfrak{p} . Returns an implicit result stream (i.e., an iterator that mutates x's variable bindings and emits rule bodies).

```
1. \operatorname{def} \mathfrak{p}[x]:
2. 
ightharpoonup Implicitly unify x with the heads of the rules in <math>\mathfrak{p}; return a stream of their bodies subject to implicit unification. If x is a value, return it. If x is a built-in relation, dispatch to the built-in result stream method.
3. \operatorname{if} x \in \mathbb{V}: \operatorname{yield} x
4. \operatorname{else} \operatorname{if} x is a builtin:
5. 
ightharpoonup Dispatch to builtin-specific code; example code is given in Algorithm 3.
6. <math>(\ldots)
7. \operatorname{else}
8. \operatorname{for} \mathbf{R} \in \mathfrak{p}: 
ho Use rules in the program<sup>28</sup>
9. 
ho Below, we freshen to avoid variable-name collisions in recursion
10. (x' \oplus = y_1 \circ \cdots \circ y_K) \leftarrow \operatorname{fresh}(\mathbf{R}) 
ho Often, the body of a rule is just a constant.
11. \operatorname{for}_{-} \in [x' = x]:
12. \operatorname{yield} y_1 \circ \cdots \circ y_K
```

might look like, consider the case where the query x matches the built-constraint A is B + C. The dispatching procedure looks like the pseudocode in Algorithm 3.

Notice that the pseudocode has several special cases that depend on whether or not A, B, and C are bound to actual numbers. In the final "else" case, we see an issue regarding built-in relations: they must be *sufficiently instantiated* to be evaluated. For example, the addition builtin A is B+1 is not sufficiently instantiated: it refuses to enumerate the infinitely many pairs $\langle A, B \rangle$ that satisfy the relation. Determining when a relation is sufficiently instantiated is left to the built-in library designer's discretion. Thus, by modifying Algorithm 3, the library *could* (lazily) enumerate all infinitely many pairs of floating-point numbers as a proxy for the uncountable set of real solutions. Generally, built-in relations will throw instantiation fault exceptions when the number of solutions is infinite because it is more likely that the user made a mistake than they wanted an

Algorithm 3 Built-in relation dispatching.

```
match x:
2
      case A is B + C:
        if isnum(B) and B = 0:
                                                         ▷ Special cases for identity element (i.e., 0)
           yield from [A=C]
        else if isnum(c) and c = 0:
           yield from [A=B]
        else if isnum(B) and isnum(C):
                                                                          \triangleright Solve for A given B and C
           a \leftarrow \mathtt{B} + \mathtt{C}
                                               \triangleright The + operation here is part of the host language.
           yield from [A=a]
10.
        else if isnum(A) and isnum(C):
                                                                          \triangleright Solve for B given A and C
           b \leftarrow A - C
           yield from [B=b]
13.
        else if isnum(A) and isnum(B):
                                                                          \triangleright Solve for C given A and B
14
           c \leftarrow \mathsf{A} - \mathsf{B}
15.
           yield from [c=c]
16.
        else
           if optional:
                                                                               \triangleright if delaying is allowed
18.
              yield x
                                                                     ▷ returns a "delayed constraint"
           else
20
              raise InstantiationFault()
21
22
     case A is B * C:
23
24.
     case A is B / C:
26
```

infinite stream of solutions. An alternative strategy (labeled as "optional" in the code) is to *delay* the processing of the constraint in hopes that it becomes more instantiated later. We note that the implementation of certain builtins may throw runtime exceptions, and other implementations may encode errors as special values (akin to NaN in floating-point numbers). This choice is up to the designer of the specific built-in relation.

Join. Next, we describe how to *lift* the \circ operator to join streams of values and their associated variable bindings. Suppose $\beta_1 \leftarrow \mathfrak{p}[y_1]$ and $\beta_2 \leftarrow \mathfrak{p}[y_2]$ for

queries y_1 and y_2 . As suggested by the choice of notation $\beta_1 \circ \beta_2$, the meaning of joining result streams is the same as the meaning of the right-hand side of a rule with $y_1 \circ y_2$. An algorithm for creating a result stream for $\beta_1 \circ \beta_2$ is given in Algorithm 4. The algorithm is shockingly simple: it is just a pair of nested loops. It is so simple due to the implicit unification strategy: variables are seamlessly and efficiently propagated between subgoals.²⁹ In particular, there is no need to pass around or return explicit substitution mappings (θ) . These nested loops correspond to a join operation because the bindings from the outer loop over β_1 are passed to the inner loop over β_2 . The inner loop iterator is re-initialized on each iteration of the outer loop: β_2 starts from the beginning, given the initial bindings, which are (generally) different from the previous iteration of the outer loop. In other words, to ensure the correct join behavior, it is important that the iterators for $\mathfrak{p}[y_1]$ and $\mathfrak{p}[y_2]$ are only *initialized* when they are put in an iteration context (e.g., such as line 4.3 or line 4.4). This ensures that the inner loop $\mathfrak{p}[y_2]$ corresponds to a new iterator for each round of the outer loop.

Algorithm 4 $\beta_1 \circ \beta_2$: Algorithm for joining result streams β_1 and β_2 .

```
    def β<sub>1</sub> ∘ β<sub>2</sub>:
    Each iteration of the outer iterator β<sub>1</sub> binds variables through implicit unification; those bindings are propagated to the inner iterator β<sub>2</sub>.
    for v<sub>1</sub> ∈ β<sub>1</sub>:
    for v<sub>2</sub> ∈ β<sub>2</sub>:
    yield v<sub>1</sub> ∘ v<sub>2</sub>
```

The result of the join operation is another implicit result stream. We write $\mathfrak{p}[y_1 \circ \cdots \circ y_K]$ as a shorthand for $\mathfrak{p}[y_1] \circ \cdots \circ \mathfrak{p}[y_K]$. This shorthand can be formalized

²⁹Such propagation is often called a *sideways information passing strategy* in the logic programming literature.

as adding a case to Algorithm 2 that will dispatch to Algorithm 4.³⁰ Note that the for-loops in the algorithm can be performed in any order; however, the choice of order may affect the running time.³¹ Since the join operation is associative, a multi-way join need not be parenthesized; it is well-defined as the o-product of multiple result streams. In §6.3, we consider an automated method for determining which ordering of a multi-way product is most efficient to use.

3.1.2.3 Update $(\mathfrak{p}[x] \oplus = ...)$

This section discusses methods for aggregating result streams as updates to a program. Let be $\mathfrak p$ a program. Since $\mathfrak p$ is just a specification of a valuation mapping, it is a glorified map data structure (like a Python dictionary or Java hash map). Thus, it seems reasonable that a user or algorithm should be able to update a given item's value by a specific value. Therefore, we allow $\mathfrak p$ to be updated using the same notation for updating maps:³²

$$\mathfrak{p}[f(1)] \oplus = 3 \quad \stackrel{\text{translates to}}{\Longrightarrow} \quad \mathfrak{p} \leftarrow \mathfrak{p} \cup (f(1) \oplus = 3)$$

Similarly,

$$\mathfrak{p}[f(X)] \oplus = 3 \quad \overset{\text{translates to}}{\Longrightarrow} \quad \mathfrak{p} \leftarrow \mathfrak{p} \cup (f(X) \oplus = 3)$$

We can also add a new rule with subgoals

$$\mathfrak{p}[f(X)] \oplus = g(X,Y) \circ h(Y) \quad \overset{\text{translates to}}{\Longrightarrow} \quad \mathfrak{p} \leftarrow \mathfrak{p} \cup \big(f(X) \oplus = g(X,Y) \circ h(Y)\big)$$

³⁰ Another choice is to make o as builtin.

 $^{^{31}}$ Some care is needed to ensure that the multiplicands of noncommutative \circ operations appear in the correct order

³²Note that the updates here are additive (as suggested by \oplus =); they are not replacements or overrides.

Now, what happens with new rules containing subgoals that are result streams (possibly from other programs)?

$$\mathfrak{p}[f(X)] \oplus = g(X,Y) \circ \mathfrak{m}[h(Y)] \quad \overset{\operatorname{translates} \ to}{\Longrightarrow} \quad \text{for} \ v \in \mathfrak{m}[h(X)] : \quad \mathfrak{p}[f(X)] \oplus = g(X,Y) \circ v$$

Notice that the elements of the stream are not required to be values in \mathbb{V} ; they may contain delayed subgoals, for example, if \mathfrak{m} has the rule $h(X) \oplus = p(X,Z)$ in the example above. Using this notation requires some care to ensure that the meaning of the p(X,Z) items is consistent across \mathfrak{m} and \mathfrak{p} . A safe setting is one where the delayed subgoals are guaranteed to have the same meaning across programs; builtins are a good example of relations that have the same meaning across programs. This flexibility allows us to include h(X) := 3 for (1 is X mod 2) as a rule in \mathfrak{m} —broadening the set of valuations that we can compactly represent.

Lastly, we mention an important special case we implicitly handle in this method. Replacing a subgoal with a set of values (or rule bodies) makes it possible to "lose" a variable. Consider the following,

Example 7. Suppose we have an empty program \mathfrak{p} and

$$\mathfrak{m} = \begin{cases} 50 & f(X) \neq 4. \end{cases}$$

We wish to apply the rule goal += f(X) to the values in m using the following:

$$\mathfrak{p}[goal] \oplus = \mathfrak{m}[f(X)]$$

If we are not careful, the rule added to p will be goal += 4, but that would be incorrect, as the value 4 needs to be summed infinitely many times. We fix this edge case by including a 4 with an infinite multiplicity factor into the update: the

case of the real semiring $\infty \cdot 4 = \infty$. Thus, the correct update is goal += ∞ . 33

The dropped variable check is encapsulated in the process of adding a rule. In our implementation, we do the following. Let $\mathfrak{p}[x] \oplus \beta$ be an update operation where β is a stream. Let v be the current value in β 's stream. We use the following **multiplicity correction**:

$$m = \begin{cases} 1 & \text{if } \operatorname{vars}(\boldsymbol{x} \oplus \boldsymbol{\beta}) \subseteq \operatorname{vars}(\boldsymbol{x} \oplus \boldsymbol{\gamma}) \\ \infty & \text{otherwise} \end{cases}$$
 (3.1)

where $\operatorname{vars}(\cdot)$ returns the set of free variables in its argument. Here is what the condition says: if replacing the rule $x \oplus = \beta$ by $x \oplus = v$ loses a variable, we need to sum infinite copies of v. For each iterate of the stream β , we use $v \cdot m$ where m is the result of (3.1) applied to the iterate (in its current state). This correction is applied whenever we add rules using the update notation $(\mathfrak{p}[x] \oplus = \beta)$.

3.2 General Mathematics

For completeness, this section includes precise definitions of the key mathematical concepts used throughout this work.

Definition 1. A binary relation \mathcal{R} on a set \mathcal{A} is a subset of the Cartesian product: $\mathcal{R} \subseteq \mathcal{A} \times \mathcal{A}$. We typically denote inclusion tests using infix operators, e.g., $a\mathcal{R}b \iff \langle a,b \rangle \in \mathcal{R}$.

Adjectives used to define binary relations include:

³³In the case of an idempotent semiring (e.g., boolean, min-plus), the multiplicity ≥ 1 are all the same. This is why the added complexity of checking update multiplicity does not appear in the ordinary logic programming setting: the boolean semiring is idempotent.

- **Reflexive**: $\forall a \in \mathcal{A} : a \mathcal{R} a$
- Symmetric: $\forall a, b \in \mathcal{A} : a \mathcal{R} b \Longrightarrow b \mathcal{R} a$.
- **Asymmetric**: $\forall a, b \in \mathcal{A} : a \mathcal{R} b \Longrightarrow \neg (b \mathcal{R} a)$
- Antisymmetric: $\forall a, b \in \mathcal{A} : a \mathcal{R} b, b \mathcal{R} a \Longrightarrow a = b$
- Transitive: $\forall a, b, c \in \mathcal{A} : a \mathcal{R} b, b \mathcal{R} c \Longrightarrow a \mathcal{R} c$

Definition 2. A partition σ of a set \mathcal{A} is a set of (non-empty) subsets $\{\sigma_1, \dots, \sigma_K\}$ of \mathcal{A} such that each element $a \in \mathcal{A}$ is a member of precisely one of these subsets. Put differently, σ must be

- (i) mutually exclusive: $\forall i \neq j : \sigma_i \cap \sigma_j = \emptyset$
- (ii) exhaustive: $A = \bigcup_{k=1}^{K} \sigma_k$

3.2.1 Ordering Relations

Definition 3. A binary relation \equiv on a set \mathcal{A} defines an **equivalence relation** on the elements of \mathcal{A} if and only if the following conditions are satisfied for all $a, b, c \in \mathcal{A}$:

- (i) Reflexive: $(a \equiv a)$.
- (ii) Symmetric: $(a \equiv b) \Longrightarrow (b \equiv a)$.
- (iii) Transitive: $(a \equiv b)$, $(b \equiv c) \Longrightarrow (a \equiv c)$.

For a given element $a \in \mathcal{A}$ its **equivalence class** $[a]^{(\equiv)} \stackrel{\text{def}}{=} \{a' \mid a \equiv a', a' \in \mathcal{A}\}$. The set of all equivalence classes $\{[a]^{(\equiv)} \mid a \in \mathcal{A}\}$ forms a partition of \mathcal{A} .

Definition 4. A binary relation \leq on a set A defines a **partial order** on the elements of A if and only if the following conditions are satisfied for all $a, b, c \in A$:

- (i) Reflexive: $a \leq a$
- (ii) Antisymmetric: $a \leq b$, $b \leq a \Longrightarrow a = b$
- (iii) Transitive: $a \leq b$, $b \leq c \Longrightarrow a \leq c$

Definition 5. A binary relation \prec on a set \mathcal{A} defines a **strict partial order** on the elements of \mathcal{A} if and only if the following conditions are satisfied for all $a, b, c \in \mathcal{A}$:

- (i) Asymmetric: $a \prec b \Longrightarrow \neg (b \prec a)$
- (ii) Transitive: $a \prec b$, $b \prec c \Longrightarrow a \preccurlyeq c$

Definition 6. A binary relation \leq on a set \mathcal{A} defines a **total order** if and only if it is both a partial order, and it is total: $a \leq b$ or $b \leq a$ for all $a, b, c \in \mathcal{A}$.

Definition 7. A binary relation \lesssim on a set \mathcal{A} defines a **preorder** on the elements of \mathcal{A} if and only if the following conditions are satisfied for all $a, b, c \in \mathcal{A}$:

- (i) Reflexive: $a \lesssim a$
- (ii) Transitive: $a \lesssim b$, $b \lesssim c \Longrightarrow a \lesssim c$

Definition 8. A binary relation \prec over a set \mathcal{A} is a **well-founded ordering** if and only if it does not contain any **infinitely descending chains**: there is no infinite sequence $a_0, a_1, a_2, ... \in \mathcal{A}$ such that $a_{n+1} \prec a_n$ for all $n \in \mathbb{N}$.

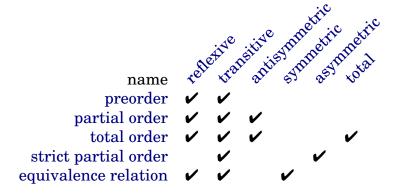


Figure 3.1: Taxonomy of order relations.

3.2.2 Functions

Definition 9. A function ϕ from a set \mathcal{A} to a set \mathcal{B} is a binary relation that assigns exactly one element of \mathcal{B} to each element of \mathcal{A} . We write $\phi : \mathcal{A} \to \mathcal{B}$ for conciseness. We often refer to $\mathcal{A} \to \mathcal{B}$ as ϕ 's type.

Definition 10. The mapping $\phi : A \to B$ is a **bijection** if for all $a \in A$:

- *surjective*: $\forall b \in \mathcal{B}$: $\exists a \in \mathcal{A}$: $\phi(a) = b$.
- *injective*: $\forall a, a' \in \mathcal{A} : (\phi(a) = \phi(a')) \Longrightarrow (a = a').$

An equivalent definition is that ϕ is a bijection from $\mathcal{A} \to \mathcal{B}$ if and only if both of the following conditions hold:

- $\forall a \in \mathcal{A} : \exists! b : \phi(a) = b$
- $\forall b \in \mathcal{B} : \exists ! a : \phi(a) = b$

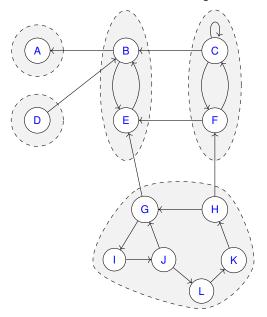
Furthermore, we say that the function ϕ^{-1} is the **inverse** of ϕ if and only if $\forall b : \phi^{-1}(b) = a$ where a is the unique solution to $\phi(a) = b$.

3.2.3 Graphs

Definition 11. A directed subgraph is **strongly connected** if and only if for each pair of nodes i and j in the subgraph, there exist paths within the subgraph from $i \leadsto j$ and $j \leadsto i$.

Definition 12. The strongly connected components (SCCs) of a graph G partition it into strongly connected subgraphs. The SCCs are the equivalence classes given by co-accessibility equivalence relation $(i \iff j) \iff (i \rightsquigarrow j)$ and $(j \rightsquigarrow i)$. The graph relating SCCs is called the **condensation** of the graph G.

Example 8. Below is a diagram³⁴ of a directed graph where its strongly connected components have been annotated with dashed ellipses.



Definition 13. A topological sort (toposort, for short) of a directed graph is a total ordering (Definition 6) of its nodes such that for every edge $i \rightarrow j$, i comes before j in the ordering. A toposort exists if and only if the graph is acyclic.

This diagram is based on https://tex.stackexchange.com/questions/592539/how-to-create-figures-with-strongly-connected-components.

We often talk about the toposort of the condensation of a cyclic graph; the condensed graph is always acyclic. For completeness, we provide pseudocode for Tarjan (1972)'s algorithm for finding a topologically sorted list of the strongly connected components of a graph in Algorithm 5.

Algorithm 5 sccs(incoming, targets): Enumerate the strongly connected components (SCC; Definition 12) of a directed graph in topological order.

1. $raiset Here the graph \ \langle V, E \rangle$ is represented implicitly using a function incoming(v) that enumerates the nodes u such that $(u \to v) \in E$, and a subset of target nodes targets $\subseteq V$ to begin traversal from; if none of the nodes in targets depend (transitively) on a given SCCs, it will be omitted from the list of SCCs returned. The algorithm is due to Tarjan (1972). It runs in $\mathcal{O}(|V| + |E|)$ time and uses $\mathcal{O}(|V|)$ additional space. Note that sccs(outgoing, V), where outgoing(u) enumerates the nodes v such that $(u \to v) \in E$, will return the same set of SCCs sccs(incoming, V), but in a reverse topological order.

```
2. def sccs(incoming, targets):
      stack ← []
                                                                  > We maintain a stack of distinct items
      pos \leftarrow map()
                                                     ▶ We record the position of each node on the stack
      \mathsf{visited} \leftarrow \{\}
                                                   ▶ We maintain a set of nodes that have been visited
      sccs_list \leftarrow []

▷ Will be the return value

      def traverse(v):
         low ← | stack |
                                                                                  ▶ Will be the return value
8
         if v \notin \mathsf{visited}:

⊳ first visit

            visited.add(v)
10
            pos[v] \leftarrow | stack |
11
            stack.push(v)
                                                                                         ▷ we may undo this
12
            for u \in \text{incoming}(v):
13
                if u \in pos:
                                                                                              ▷ cycle detected
14
                   low min= pos[u]
                else
16.
                   low min= traverse(u)
            if low = pos[v]:
                                                \triangleright If low is unchanged; nothing beneath v is in its SCC
18
                \sigma \leftarrow \{\}
19
                while | stack | > low:
20
                   u \leftarrow \mathsf{stack.pop}()
21
                   \sigma.add(u)
22.
                   pos.remove(u)
23
                sccs list.append(\sigma)
                                                                                           \triangleright v's SCC is ready
24
         return low
25
      for v \in \text{targets}: traverse(v)
26
      return sccs list
```

3.2.4 Abstract Algebra

Definition 14. A binary operation on a set \mathbb{V} is any function $\odot : \mathbb{V} \times \mathbb{V} \to \mathbb{V}$. Binary operators are often written using infix notation $a \odot b$ rather than $\odot (a, b)$, which is typical for function applications.

Definition 15. A monoid $\langle \mathbb{V}, \odot, \underline{i} \rangle$ is a set \mathbb{V} equipped with a binary operation \odot and an identity element \underline{i} such that the following properties hold:

- (i) \odot is associative: $\forall a, b, c \in \mathbb{V} : (a \odot b) \odot c = a \odot (b \odot c)$
- (ii) \underline{i} is the left and right identity element: $\forall a \in \mathbb{V} : \underline{i} \odot a = a \odot \underline{i} = a$

Some adjectives used to describe monoids include:

- **commutative** $\iff \forall a, b \in \mathbb{V} : a \odot b = b \odot a$
- idempotent $\iff \forall a \in \mathbb{V} : a \odot a = a$

Definition 16. A group is a monoid $\langle \mathbb{V}, \odot, \underline{i} \rangle$ where additionally, every element has an inverse element: $\forall a \in \mathbb{V}, \exists b \in \mathbb{V} : b \odot a = a \odot b = \underline{i}$. The element b may be denoted as a^{-1} if the group is considered product-like or -a if the group is considered sum-like.

Definition 17. A semiring $\langle \mathbb{V}, \oplus, \circ, \underline{0}, \underline{1} \rangle$ is a set \mathbb{V} equipped with two monoid operations $\langle \mathbb{V}, \oplus, \underline{0} \rangle$ and $\langle \mathbb{V}, \circ, \underline{1} \rangle$ such that the following properties hold:

- (i) $\langle \mathbb{V}, \oplus, \underline{0} \rangle$ is a commutative monoid
- (ii) $\langle \mathbb{V}, \circ, \underline{1} \rangle$ is a monoid

(iii) Distributivity: $\forall a, b, c \in \mathbb{V}$:

$$a \circ (b \oplus c) = (a \circ b) \oplus (a \circ c)$$
 and $(b \circ a) \oplus (c \circ a) = (b \oplus c) \circ a$

(*iv*) Multiplicative annihilation by $\underline{0}$: $\forall a \in \mathbb{V}$:

$$0 \circ a = a \circ 0 = 0$$

Additionally, we adopt the following notation for the n-fold sum

$$b \cdot n \stackrel{\text{def}}{=} n \cdot b \stackrel{\text{def}}{=} \bigoplus_{i=1}^{n} b \quad \text{for } n \in \mathbb{N} \cup \{\infty\}, b \in \mathbb{V}$$
 (3.2)

and the *n*-fold product

$$b^{n} \stackrel{\text{def}}{=} \bigotimes_{i=1}^{n} b \quad \text{for } n \in \mathbb{N} \cup \{\infty\}, b \in \mathbb{V}$$
 (3.3)

Some adjectives used to describe semirings include:

- commutative ⇔ ∘ is commutative
- idempotent \iff \oplus is idempotent

Definition 18. A closed semiring has an additional operator, called the (Kleene) star operator $(\cdot)^*$. This operation must satisfy the following axioms for all $a \in \mathbb{V}$:

$$a^{\star} = 1 \oplus a \circ a^{\star}$$
 and $a^{\star} = 1 \oplus a^{\star} \circ a$

We often interpret a^* as the solution to the infinite series that arises from unrolling either of its recursive identities:

$$a^\star = igoplus_{n=0}^\infty a^n \quad ext{for } a \in \mathbb{V}$$

The star operation is useful for solving certain linear equations. For example, $x = b \oplus a \circ x \oplus c \circ y$ can be rearranged to "solve" for x without subtraction or division: $x = a^* \circ (b \oplus c \circ y)$.

A special family of semirings that always have a star operation is those that are *K*-closed.

Definition 19. For $K \in \mathbb{N}$, a semiring $(\mathbb{V}, \oplus, \circ, \underline{0}, \underline{1})$ is K-closed if and only if 85

$$\forall \, a \in \mathbb{V}, \forall \, K' > K \colon \bigoplus_{k=0}^{K'} a^k = \bigoplus_{k=0}^{K} a^k.$$

Intuitively, K-closed semirings imply that we only have to sum contributions of a cycle by iterating around the cycle at most K times; further iterations do not affect the value of the summation. When K=0, the condition simplifies to $\forall a \in \mathbb{V} : \underline{1} \oplus a = \underline{1}$. We briefly mention some examples of K-closed semirings; these examples are defined in §3.3. Examples of 0-closed semirings include the boolean, min-plus, and Viterbi semirings. An example of a general K-closed semiring include the (K+1)-best semirings (Mohri, 2002, §6) and the K-collapsed natural numbers. Non-examples of K-closed semirings include the arithmetic semirings.

Each K-closed semiring has simple expressions for computing the star operation because they are no longer infinite sums.

 $[\]overline{^{35}}$ It is possible to replace $\forall K'$ with $\overline{K'}=K+1$ —the conditions are equivalent (Mohri, 2002, Lemma 4).

$$\begin{cases} 0\text{-closed}: a^* = \underline{1} \\ 1\text{-closed}: a^* = \underline{1} \oplus a \\ 2\text{-closed}: a^* = \underline{1} \oplus a \oplus a^2 \\ \vdots \\ K\text{-closed}: a^* = \bigoplus_{k=0}^K a^k \end{cases}$$

$$(3.4)$$

Definition 20. A semiring is **monotonic** under the partial ordering relation \leq if

$$\forall a \preccurlyeq b, \forall c \in \mathbb{V} : \begin{cases} a \oplus c \preccurlyeq b \oplus c \\ a \circ c \preccurlyeq b \circ c \\ c \circ a \preccurlyeq c \circ b \end{cases}$$

Where to now? In §3.3, we will describe a number of examples of semirings. In §4.4, we will describe some special families of semirings that are useful in providing convergence guarantees for Dyna programs.

3.3 Semirings: A Powerful Abstraction for Dynamic Programming

Semirings are a powerful abstraction for dynamic programming. They provide a useful and elegant abstraction level for

- algorithms for NLP (e.g., Teitelbaum (1973), Mohri (1997), Goodman (1999),
 and Huang (2008))
- probabilistic inference (e.g., Aji and McEliece (2006))
- database provenance (e.g., Green, Karvounarakis, et al. (2007), Green and Tannen (2017), and Cheney et al. (2009))
- network analysis (e.g., Gondran and Minoux (2008) and Baras and Theodorakopoulos (2010))
- program analysis (e.g., Cousot and Cousot (1977), Esparza et al. (2010),
 and Dolan (2013))

Semirings enable the following separation of concerns:

- (1) devising a compact encoding of a set of derivations (e.g., paths in a graph, parses of a sentence)
- (2) how to evaluate and aggregate the derivations (e.g., the highest-weight derivation, the total weight of all derivations)

In Dyna, part 1 is specified by the program rules, and part 2 is given by choice of semiring. Modularizing by semiring gives different semantics for the same algorithm with the same time/space complexity (assuming semiring operations have unit time), and we can even reuse the implementation.³⁶

The key property of a semiring is distributivity, as it allows us to factor out common subexpressions: we can safely rewrite $(a \circ b) \oplus (a \circ c)$ (three operations) into an equivalent smaller expression $a \circ (b \oplus c)$ (two operations). The factored expression can be significantly smaller.

Example 9. The example below may be regarded as encoding a weighted directed graph with two states \times and y. The weights in the graph have been left abstract as symbols a, b, c, and d.

$$\begin{cases} x += a * x. \\ x += b * y. \\ y += c * x. \\ y += d * y. \\ y += 1. \% \text{ final state} \end{cases}$$

We can translate the Dyna program into a system of equations:

$$x = \underset{\bullet}{\mathbf{a}} \circ x \oplus \underset{\bullet}{\mathbf{b}} \circ y$$
$$y = \underset{\bullet}{\mathbf{c}} \circ x \oplus \underset{\bullet}{\mathbf{d}} \circ y \oplus \underline{1}$$

The solution to the system of equations is a map where the value associated with x is the set of all paths starting in state x and ending in state y, and the value associated with y is the set of all paths starting in state y and ending in state y. In this graph, there are infinitely many paths, and yet, we can express them as a small set of equations. Below we consider replacing the weights (a, b, c, d) with

³⁶Although certain semirings may allow for additional efficiency shortcuts that others do not, certain implementations may assume properties beyond the basic semiring axioms.

values in a variety of semirings. In each case, we describe the interpretation of the semiring values associated with \times and y.

- $\langle \{ \text{false, true} \}$, \vee , \wedge , $\{ \text{false, true} \}$ where true is placed on each edge \Longrightarrow y is reachable from x, and y is reachable from y.
- $\bullet \ \left\langle \mathbb{R} \cup \left\{ \infty \right\}, +, \cdot, \infty, 0 \right\rangle \Longrightarrow \textit{the total weight of all paths from x to y and y to y}.$
- $\langle \mathbb{R}, \min, +, 0, 1 \rangle \Longrightarrow \textit{the least cost path from} \times \textit{to} \ \mathsf{y} \ \textit{and} \ \mathsf{y} \ \textit{to} \ \mathsf{y}.$
- $\langle 2^{\Sigma^*}, \cup, \cdot, \emptyset, \varepsilon \rangle$ where Σ is the alphabet $\{a, b, c, d\} \Longrightarrow$ the (regular) formal language that describes the set of paths from x to y and from y to y.

This is just a small example of what we can do by swapping the semiring and keeping the Dyna program essentially the same.

Semiring Compendium

Below we list several commonly used semirings and some applications where they appear. Many of these semirings are available in the Dyna library https://github.com/timvieira/dyna-pi. We do not have a particular organizing principle for this list, but we hope that it may serve as a useful reference for some readers and that it illustrates that semirings can do some surprising things.

Logic

Boolean semiring $\langle \{false, true\}, \vee, \wedge, false, true \rangle$ are used to test reachability in graphs and, more generally, provability of items under the rules of the

program. Its star operation is $(false)^* = (true)^* = true$

Multi-valued logic.

- Three-valued logic $\{ \text{true}, \text{false}, \text{unk} \}$, $\max, \min, \text{false}, \text{true} \}$ where false < unk < true. 37
- Fuzzy sets $\langle [0,1], \max, \min, 0, 1 \rangle$ is related to fuzzy sets (Zadeh, 1965).³⁸
- Łukasiewicz semiring $\langle [0,1], \max, \max(0, x+y-1), 0, 1 \rangle$ (Łukasiewicz, 1920)

These semirings are each 0-closed; thus, their star operations can be obtained using the K-closed shortcut (3.4).

Arithmetic semirings

- Counting semiring: $\langle \mathbb{N}, +, \cdot, 0, 1 \rangle$. As the name suggests, this semiring is used to count things, such as the number of paths between two nodes in a graph. If \mathbb{N} is augmented with ∞ , then $a^* = 1$ if a = 0 else ∞ .
- Integer semiring: $\langle \mathbb{Z}, +, \cdot, 0, 1 \rangle$. Its star operation is $0^* = 1$, $a^* = \infty$ if a > 0, and undefined if a < 0.
- **Real semiring**: $\langle \mathbb{R}, +, \cdot, 0, 1 \rangle$. If we extend \mathbb{R} with ∞ , we have $a^* = \frac{1}{1-a}$ if |a| < 1, ∞ if $a \ge 1$, and it is undefined if $a \le -1$. So, unfortunately, this semiring does not have complete support for solving the infinite series.

³⁷Often \langle false, unk, true \rangle are encoded as $\langle -1, 0, +1 \rangle$.

³⁸A fuzzy set extends an ordinary set to support graded membership that indicates the extent to which a given element belongs to a given set (0 meaning definitely not in the set, and 1 meaning definitely in the set).

- Nonnegative-real semiring: $\langle \mathbb{R}_{\geq 0} \cup \{\infty\}, +, \cdot, 0, 1 \rangle$. This semiring is used for computing the normalization constant of a probabilistic model (Goodman, 1999). Its star operation is $a^* = \frac{1}{1-a}$ for a < 1 and ∞ otherwise. Thus, to be closed, we require ∞ to be added to the set \mathbb{R} .
- Log-nonnegative-real semiring: $\left\langle \mathbb{R} \cup \{-\infty,\infty\}, +, +, -\infty, 0 \right\rangle$ where $x + y \stackrel{\text{def}}{=} \log(\exp(x) + \exp(y))$. This semiring is isomorphic to the nonnegative-real semiring; however, it often preferred as it can avoid certain numerical issues. Its star operation is $a^* = -\log(1 \exp(a))$ if $a \leq 0$ and ∞ otherwise.
- **First-order expectation semirings** (Eisner, 2001) can be used to compute first-order statistics, such as the expected error rate, the entropy, or the expected feature counts.
- **Second-order expectation semirings** (Li and Eisner, 2009) can be used to compute second-order statistics, such as the variance of the error rate, the gradient of entropy, and the covariance matrix of feature counts.
- K-collapsed natural numbers $\langle \{0,\cdots,K\}, \min(K,x+y), \min(K,x\cdot y), 0, 1 \rangle$ can be used as guard rails that prevent infinite running times due to counting off to infinity. It is isomorphic to the boolean semiring when K=2. Its star operation is $a^*=1$ if a=0 else K.

³⁹For speed and numerical precision, the log-space addition operation should not be computed as written, it should be computed using a library function (Mächler, 2012). The numpy library for Python provides the function numpy logaddexp for this purpose.

⁴⁰We also mention the log-real semiring (Li and Eisner, 2009, Table 3), which extends the log-nonnegative-real semiring with an additional sign bit so that negative numbers can be represented.

Optimization

These semirings are often used when computing the minimum/maximum values over some set of values, for example, the weighted set of paths in a graph.

- The **tropical semiring** (aka **min-plus semiring**) $\langle \mathbb{R} \cup \{\infty\}, \min, +, \infty, 0 \rangle$ is used to find the shortest path or cheapest derivation. Its star operation is $a^* = 0$
- The arctic semiring (aka max-plus semiring) $\langle \mathbb{R} \cup \{-\infty\}, \max, +, -\infty, 0 \rangle$ is used to find the longest path or highest-weight derivation. Its star operation is $a^* = \underline{1} \oplus a$ if $a \leq 0$ else ∞
- The **Viterbi semiring** (aka **max-times semiring**) $\langle [0,1], \max, \cdot, 0, 1 \rangle$ is used to find the most probable explanation (MPE) in a probabilistic model. MPE corresponds to maximum a posteriori (MAP) decoding when the model does not have latent variables. Its star operation is $a^* = 1$. If we broaden the carrier set with numbers > 1, we have to extend the definition to $a^* = \infty$ for a > 1.
- The **min-times semiring** $\langle \mathbb{R}_{\geq 0} \cup \{\infty\}$, $\min, \cdot, \infty, 1 \rangle$ can be used to find the least probable derivation. Its star operation is $a^* = \underline{1} \oplus a$

These semirings can be augmented with additional bookkeeping to return an argmin/argmax (or even the set of argmins/argmaxes) in addition to the min/max. It is also possible to make the elements being compared in these optimization

semirings tuples and to compare them under a lexicographic ordering. This is one way to keep track of back pointers and to add tie-breaking operations.

The **maximum-capacity semiring** $\langle \mathbb{R} \cup \{-\infty, \infty\}$, $\max, \min, -\infty, \infty \rangle$ can be used to find the maximum capacity path in a graph. It is 0-closed so $a^* = \underline{1}$.

The K-best semiring (top-K semiring or bottom-K semiring) extends the optimization semirings to provide the K-best alternatives where $K \geq 0$ is a fixed constant.⁴¹ The K-best semirings work with any monotonic product operations, which include multiplication (by nonnegative numbers), Cartesian products, and addition. Let $\langle \mathcal{A}, \circ, \varepsilon \rangle$ be a monoid. Suppose that it is partially ordered under \leq and the \circ operation is monotonic.

- $\mathbb{V} \stackrel{\text{def}}{=} \{ \operatorname{sort}_{\preccurlyeq}^K(a_1 \cdots a_K) \mid a_1 \dots a_K \in \mathcal{A}^K \}$: each element of this set is a length-K sorted list of elements from \mathcal{A} . Here the function $\operatorname{sort}_{\preccurlyeq}^K$ takes a sequence, sorts it under \preccurlyeq , and truncates only to include the K smallest elements.
- $\bullet \ A \oplus B \stackrel{\mathrm{def}}{=} \mathrm{sort}_{\preccurlyeq}^K (A \cup B)$
- $\bullet \ A \circ B \stackrel{\mathrm{def}}{=} \mathrm{sort}_{\preccurlyeq}^K([\underline{a} \circ \underline{b} \mid \underline{a} \in A, \underline{b} \in B])$
- 0 def []
- $\underline{\underline{1}} \stackrel{\text{def}}{=} [\varepsilon]$

 $[\]overline{^{41}}$ Goodman (1998) gave the *K*-best semiring, which eagerly combines the top-*K* values for some fixed *K*. Later, this approach was generalized (Huang and Chiang, 2005) to allow the *K*-best list to be computed lazily (and without knowledge of *K*), which can be significantly faster as many of those *K*-best lists are not even used (e.g., a *K*-best list of a span that does not appear in any complete parse). The implementation details for lazily combining streams efficiently are given in Huang and Chiang (2005).

Since this semiring is (K-1)-closed. We can use the truncated sum to implement the star operation $a^* = \sum_{k=0}^{K-1} a^k$.

Matrices

The semiring of matrices (Lehmann, 1977) is important as we will use it in §5.3 to solve certain linear families of infinite sums.

Definition 21 (Matrices over a semiring). We can "lift" any base semiring $\langle \mathbb{V}, \oplus, \cdot, \underline{0}, \underline{1} \rangle$ into a semiring over square matrices $\langle \mathbb{V}^{D \times D}, +, \underline{\,\,\,\,}, 0, 1 \rangle$ where D is a set.⁴³ The matrix addition and multiplication operations (given below) are a straightforward generalization of the usual matrix multiplication on \mathbb{R} . Let $A, B \in \mathbb{V}^{D \times D}$.

- $ullet \ [\mathbf{A} + \mathbf{B}]_{ij} \stackrel{ ext{def}}{=} \mathbf{A}_{ij} \oplus \mathbf{B}_{ij}$
- $\bullet \ [\mathbf{A}\,\mathbf{B}]_{ik} \stackrel{\text{def}}{=} \bigoplus_{j \in D} \mathbf{A}_{ij} \circ \mathbf{B}_{jk}$
- $[\mathbf{0}]_{ij} \stackrel{\text{def}}{=} \underline{\mathbf{0}}$
- $[1]_{ij} \stackrel{\text{def}}{=} \underline{1}$ if i = j else $\underline{0}$

The star operation A^* can be computed as follows. If |D| = 1, then $A^* = [a]^* = [a^*]$. Otherwise, |D| > 1 and we can decompose A into nontrivial blocks $\langle B, C, D, E \rangle$ and apply the following blockwise recursion.

$$\mathbf{A}^* = \begin{bmatrix} \mathbf{B} & \mathbf{C} \\ \mathbf{D} & \mathbf{E} \end{bmatrix}^* = \begin{bmatrix} \mathbf{B}^* + \mathbf{B}^* \mathbf{C} \, \mathbf{\Delta}^* \, \mathbf{D} \, \mathbf{B}^* & \mathbf{B}^* \mathbf{C} \, \mathbf{\Delta}^* \\ \mathbf{\Delta}^* \, \mathbf{D} \, \mathbf{B}^* & \mathbf{\Delta}^* \end{bmatrix}$$
(3.5)

 $[\]overline{^{43}\text{We have denoted the juxtaposition operation (e.g., <math>\mathbf{A}\,\mathbf{B})}$ as \Box , the so-called "visible space."

where $\Delta = E + D B^*C$. It is straightforward to verify that each recursive call to $(\cdot)^*$ is decreasing in size; thus, the recursive definition is well-founded. Furthermore, if the blocks are chosen to be equal in size, then the runtime is $\mathcal{O}(|D|^3)$. We refer the reader to Lehmann (1977) for proof of correctness and further discussion.

Provenance semirings

Provenance semirings (Green, Karvounarakis, et al., 2007; Green and Tannen, 2017; Cheney et al., 2009; Schlund, 2016) are used in database systems to validate or explain query results. Similarly, "derivation semirings" (Goodman, 1998) can be used to encode the set of "derivations" for a given item. ⁴⁴ Some examples of provenance semirings given in Green, Karvounarakis, et al. (2007) include the following:

- The **set semiring** $\langle 2^{\Sigma}, \cup, \cap, \emptyset, \Sigma \rangle$ where Σ is a base set can be used to track the possibility of events in a system.
- The **why semiring** $\left\langle 2^{2^{\Sigma}}, \cup, \left\{ x \cup y \right\}_{x \in X, y \in Y}, \emptyset, \left\{ \emptyset \right\} \right\rangle$ where Σ is a base set. The base set Σ is a set of annotations; the values computed with this semiring are the combinations of annotations required for a given tuple to appear in answer to a query.
- The lineage semiring $\langle 2^{\Sigma} \cup \{\bot\}, \cup, \cup, \bot, \emptyset \rangle$ where Σ is a base set and \bot is

⁴⁴Note: Goodman's notation of a "derivation" is different from our notion (which is formalized later in § 4). In particular, his notion is less general than ours because ours is not built from semiring operations: our notion of a derivation is given by nested rules, which are not multiplicatively associative.

a special element satisfying the semiring properties of <u>0</u>. In the context of a graph, the lineage semiring determines the subset of reachable arc labels.

Formal Languages

Formal languages: $\langle 2^{\Sigma^*}, \cup, \cdot, \emptyset, \varepsilon \rangle$ where Σ is a set of symbols and Σ^* denotes the set of strings composed of those symbols, ε is the empty string, and \cdot denotes string concatenation. The semiring of formal languages can be used to enumerate the set of strings accepted by a context-free grammar, a finite-state automaton, and many other formalisms. There are, unfortunately, no shortcuts for computing the star operation. Finite-state automata can compactly represent certain infinite sets of strings; they are closed on the union, concatenation, and star operations.⁴⁵

Weighted formal languages: The semiring of formal languages can be extended to semiring-weighted languages, where a semiring weight is assigned to each string. Let $\langle \mathbb{V}, \oplus, \circ, \underline{0}, \underline{1} \rangle$ be a base semiring. Let Σ be an alphabet as in the semiring of formal languages. The semiring of weighted formal languages $\langle \widehat{\mathbb{V}}, \widehat{\oplus}, \widehat{\circ}, \widehat{0}, \widehat{1} \rangle$ is defined as follows. Let $\widehat{\mathbb{V}} = \mathbb{V}^{\Sigma^*} = \Sigma^* \to \mathbb{V}$ be the set of all functions from strings in Σ^* to values in the base semiring. Let $a, b \in \widehat{\mathbb{V}}$ and $x \in \Sigma^*$:

•
$$[a \widehat{\oplus} b](x) \stackrel{\text{def}}{=} a(x) \oplus b(x)$$

⁴⁵Although context-free languages (CFLs) are also closed under the semiring operations, the representations of CFLs (i.e., context-free grammars and pushdown automata) do not have (complete) algorithms to test the equality of the context-free languages that they represent, severely limiting their utility as a computational representation.

•
$$[a \circ b](x) \stackrel{\text{def}}{=} \bigoplus_{x: x=y \cdot z} a(y) \circ b(z)$$

$$\bullet$$
 $[\widehat{\mathbf{0}}](x) \stackrel{\text{def}}{=} \underline{\mathbf{0}}$

•
$$[\widehat{\mathbf{1}}](x) \stackrel{\text{def}}{=} \underline{\mathbf{1}}$$
 if $x = \varepsilon$ else $\underline{\mathbf{0}}$

There are, unfortunately, no shortcuts for computing the star operation.

Sets of Vectors

- The **semiring of semilinear sets** is $\langle 2^{\mathbb{N}^{\Sigma}}, \cup, +, \emptyset, \{0\} \rangle$ where $X + Y \stackrel{\text{def}}{=} \{x + y\}_{x \in X, y \in Y}$. This semiring can be used to compute the Parikh image of a formal language with alphabet Σ . We refer the reader to Schlund (2016, Ch. 5) for further discussion and efficient implementation.
- The **Pareto frontier semiring** is $\langle \text{Pareto efficient sets of vectors, pareto}(X \cup Y), \text{pareto}(X+Y), \emptyset, \{0\} \rangle$. The function pareto filters the set of vectors to those that are Pareto efficient. A set of vectors is **Pareto efficient** if and only if each of its vectors is non-dominated. A vector is said to dominate another vector if it is \leq all the other vectors in all dimensions and < in at least one dimension. The set of Pareto efficient vectors represents the optimal trade-offs between different objectives or criteria (i.e., where improving one objective comes at the cost of worsening another). The Pareto frontier semiring on $\mathbb{R}^D_{\geq 0}$ is 0-closed, so its star operation is $X^* = \underline{1}$. This semiring is discussed further in Gondran and Minoux (2008).

⁴⁶The **Parikh image** of a language L is defined as $\operatorname{Parikh}(L) \stackrel{\operatorname{def}}{=} \{ \operatorname{bag}(w) \mid w \in L \}$ where $\operatorname{bag}(w) \stackrel{\operatorname{def}}{=} \{ a \mid a \in w \}$. Placing the symbols of each word into a bag removes their order.

• The semiring of convex sets in a vector space is $\langle \text{convex sets over a vector space}, \text{convexhull}(X \cup Y), \text{convexhull}(X+Y), \emptyset, \{0\} \rangle$. For a finite-dimensional vector space, we may restrict from convex sets to polytopes and still have a semiring—the advantage being that each polytope can be computationally represented as a finite set of vertices, yielding polytime operations. Dyer (2013) showed how to use this semiring to perform efficient minimum error rate training (Och, 2003), a common technique used in machine translation. An earlier citation for this semiring is Iwano and Steiglitz (1990).

Other

- Polynomials, and formal power series (Kuich, 1997; Green, Karvounarakis,
 et al., 2007) over semirings form a semiring in each case.
- Baras and Theodorakopoulos (2010) and Gondran and Minoux (2008) give many interesting semirings for network analysis, including semirings to enumerate graph cuts and bridges.
- Many works use semirings to perform dataflow analyses on programs (Dolan, 2013; Esparza et al., 2010; Cousot and Cousot, 1977).

⁴⁷The semiring of convex sets is similar to the Pareto frontier semiring. The difference is that convex sets find the set of optimal derivations with respect to linear combinations of the base objectives.

Chapter 4

Semantics of Dyna Programs

In §2, we introduced the Dyna language through a number of examples. In this chapter, we will formalize the meaning of a Dyna program—i.e., what it actually computes—as a certain sum over potentially infinitely many semiring-weighted derivations (§4.1).⁴⁸ We illustrate the challenges with several simple examples. Quite surprisingly, we can often compute—or, at least accurately approximate—these infinite sums. We describe fixpoint iteration as the mathematical algorithm framework for computing these semantics. We later (§5) reify the framework in executable algorithms.

⁴⁸The prior formalization of Dyna (Eisner, Goldlust, et al., 2005) did not describe Dyna's semantics as a sum over weighted derivations, which meant that it was mismatched with other work on semiring-weighted deduction, e.g., Goodman (1999). This mismatch, however, enables the generalization of Dyna 1 to the Dyna 2 language (Eisner and Filardo, 2011).

4.1 Total Value of All Derivations

Example 10. Consider the following program, which encodes a geometric series:

$$\mathfrak{p} = \begin{cases} 56 & \mathsf{x} += \mathsf{b}. \\ 57 & \mathsf{x} += \mathsf{a} * \mathsf{x}. \end{cases}$$

At a high level, we generate these derivations by repeated substitution in the rules matching a subgoal.⁴⁹ Below, we show these rule substitutions by nesting rules within other rules. Each of these nested rules encodes a derivation tree.

⁴⁹We will make this formal in Definition 23.

```
58  x += b.

59  x += a * (x += b).

60  x += a * (x += a * (x += b)).

61  x += a * (x += a * (x += a * (x += b))).

62  x += a * (x += a * (x += a * (x += a * (x += b)))).

63  x += a * (x += b))))).

64  x += a * (x += b))))).

65  ...
```

Next, we can summarize each derivation by its value—i.e., the product of the values of its leaves.⁵⁰ By doing so, we recover the terms of the infinite series:

Notice that our ability to expand \times into a set of derivations whose values can be summed requires distributivity.

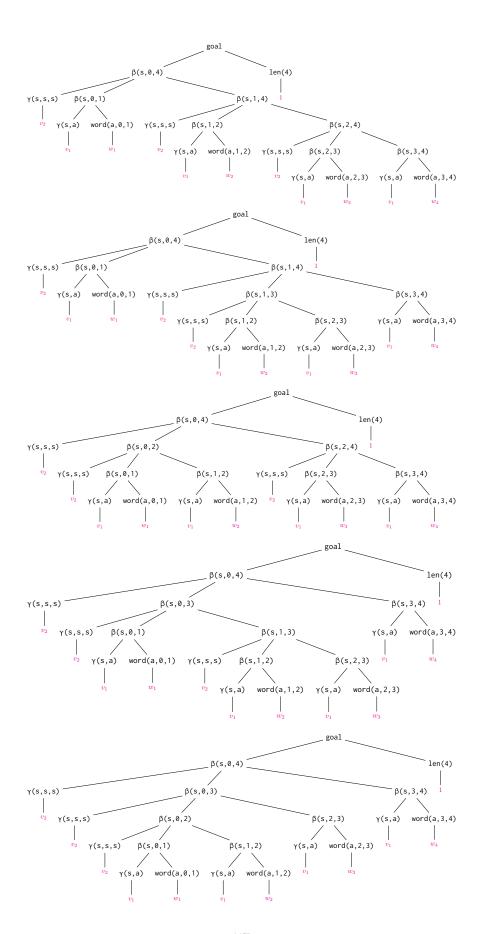
We now look at an example program containing variables.

Example 11. Consider CKY (Example 6; repeated below) where we have augmented the original program with some input data:

```
74 \beta(X,I,K) += \gamma(X,Y,Z) * \beta(Y,I,J) * \beta(Z,J,K). 78 % context-free grammar rules 
75 \beta(X,I,K) += \gamma(X,Y) * \beta(Y,I,K). 79 \gamma(s,a) += v_1. 76 \beta(X,I,K) += \gamma(X,Y) * word(Y,I,K). 80 \gamma(s,s,s) += v_2. 77 goal += \beta(s,0,N) * len(N). 81 % encoding of the input sentence "a a a a" 82 word(a,0,1) += w_1. word(a,1,2) += w_2. word(a,2,3) += w_3. word(a,3,4) += w_4. 83 len(4) += 1.
```

where $v_1, v_2, w_1, w_2, w_3, w_4 \in \mathbb{V}$. Now, consider the goal item. It has five derivations under the CKY program's rules:

⁵⁰The product uses an in-order traversal of the tree given by the nested terms to preserve the order of possibly noncommutative product operations.



Again, we see that each derivation's internal nodes are labeled with items, and its leaves are semirings values. Each derivation's value is the product of its leaves.

To determine the value of the goal we sum the values of each of its derivations:

The other items in the program also have derivations, and values derived from those derivations, but we will not show those to keep the example brief.

Next, we formalize the concepts introduced in the example above.

Definition 22. The meaning of a Dyna program⁵¹ \mathfrak{p} is a **valuation** $\mathfrak{p}[\![\cdot]\!]:\mathbb{H}\to \mathbb{V}$, ⁵² which computes the sum over derivation values. ⁵³

$$\mathfrak{p}\llbracket x \rrbracket \stackrel{\text{def}}{=} \begin{cases} x & \text{if } x \in \mathbb{V} \\ \mathfrak{B}\llbracket x \rrbracket & \text{if } x \in \mathfrak{B} \\ \bigoplus_{\delta \in \mathfrak{p}\{x\}} \mathsf{yield}(\delta) & \text{otherwise} \end{cases} \tag{4.1}$$

where $\mathfrak{B}[\![x]\!]$ is the valuation for built-in relations (items in \mathfrak{B}),⁵⁴ and $\mathfrak{p}\{x\}$ (Definition 23, below) is the (possibly infinite) set of all derivations of the item x and $yield(\delta) \in \mathbb{V}$ is the product of the values in the leaves of the derivation δ (in left-to-right order). We note that when $\mathfrak{p}\{x\}$ is infinite, (4.1) might be undefined

⁵¹Note that the Dyna $\underline{2}$ semantics (discussed in $\S D$) are *not* based on such a sum over derivations. ⁵²We may sometimes denote the set of valuations as $\mathbb{V}^{\mathbb{H}}$ rather than $\mathbb{H} \to \mathbb{V}$.

We may sometimes denote the set of variations as V rather than $\mathbb{H} \to V$.

53 Typically, $\mathfrak{p}[\cdot]$ assigns non-0 values to only a finite subset of \mathbb{H} (that are not builtins).

⁵⁴Unlike program items, built-in relations are defined directly by their valuation; In other words, the valuation for builtins is opaque: we do not have a collection of rules that define the builtins—they are essentially just foreign procedure calls. For simplicity, we assume that the set of builtins (i.e., $x \in \mathfrak{B}$) is disjoint from the items defined in the program.

(e.g., Example 12, below). When the program $\mathfrak p$ is clear from context, we may write $\llbracket \cdot \rrbracket$ as a shorthand for $\mathfrak p \llbracket \cdot \rrbracket$.

Definition 23. We define a **derivation** δ in program \mathfrak{p} as a tree structure that satisfies one of the following conditions:

- (i) $\delta \in \mathbb{V}$
- (ii) δ is of the form $x \oplus v$ where x is a builtin and $v = \mathfrak{B}[x]$
- (iii) δ is a product of derivations $\delta_1 \circ \cdots \circ \delta_K$
- (iv) δ is of the form $(\mathbf{x} \oplus = \delta_1 \circ \cdots \circ \delta_K)$ where
 - (a) δ_k are valid derivations for each k
 - (b) $x \oplus = \text{head}(\delta_1) \circ \cdots \circ \text{head}(\delta_K)$ is a grounding of a rule in $\mathfrak p$ where $\text{head}(\delta_k)$ gives δ_k 's root label for each k. 55

These cases define a labeled tree with its internal nodes labeled by rules (or built-in relations) and its leaves labeled by values. Notice that valid derivations may contain variables—meaning that any grounding of its variables is also a valid derivation. The **value** of a derivation δ is yield(δ); we sometimes write $[\![\delta]\!]$ intead of yield(δ). Finally, let $\mathfrak{p}\{x\}$ denote the **set of all derivations** of the item x. Technically, we want the smallest set of objects satisfying the description of a derivation, as this disallows derivations of infinite height. $[\![\delta]\!]$

⁵⁵Recall that we also use $head(\cdot)$ to access the head of a rule; the reuse of this notation is not coincidental, as each derivation is built from nested rules (with ground variables).

⁵⁶The set of derivations can often be represented compactly using a *packed forest* data structure, which is also known as a *B-hypergraph* (Gallo et al., 1993). By analogy, a *set of derivations* is to a *packed forest* as a *set of strings* is to a *finite-state automaton*. In both cases, we can represent certain possibly infinite sets of the respective structures efficiently by exploiting structure sharing. Packed forests are commonly used in NLP to represent parse forests and possible translation outputs (e.g., Klein and Manning (2001), Li and Eisner (2009), Huang (2008), and Eisner, Kornbluh, et al. (2006)).

Definition 24. The **height** of a derivation is defined as follows

$$\mathsf{height}(\pmb{\delta}) \stackrel{\mathsf{def}}{=} \begin{cases} 0 & \textit{if } \pmb{\delta} \textit{ is a constant or builtin} \\ \max(\mathsf{height}(\pmb{\delta}_1), \dots, \mathsf{height}(\pmb{\delta}_K)) & \textit{if } \pmb{\delta} \textit{ is a product } \pmb{\delta}_1 \circ \cdots \circ \pmb{\delta}_K \\ 1 + \mathsf{height}(\mathsf{body}(\pmb{\delta})) & \textit{if } \pmb{\delta} \textit{ is a derivation by a rule} \end{cases}$$

Not all infinite sums converge to a well-defined value. If a given item in the program does not have a well-defined sum over derivations, then its semantics are also not well-defined. In §4.4, we describe the family of ω -continuous semirings (Kuich, 1997), which are guaranteed to have a well-defined value. However, we do not require the user to use such semirings.

4.2 Fixpoint Approximation

Definition 25. We can rearrange (4.1) into a system of equations, known as the **fixpoint equations**⁵⁷ for a given program \mathfrak{p} .⁵⁸

Recall that $\Gamma(\mathbf{R})$ is the set of all groundings of the rule \mathbf{R} (defined in §2).

The fixpoint equations suggest a recursive algorithm for evaluating $\llbracket \cdot \rrbracket$. That algorithm would compute $\llbracket x \rrbracket$ by recursively computing the values of all $\llbracket y_k \rrbracket$ items that $\llbracket x \rrbracket$ depends on and plugging them into (4.2). However, this algorithm will recurse forever anytime it encounters a cycle.

⁵⁷A **fixpoint equation** has the form x = f(x) for some vector-valued function f.

⁵⁸Dyna <u>2</u> does not have the sum-over-derivations semantics—its semantics are any fixpoint of a generalized version of (4.2). For further discussion, we refer the reader to §D.

Another important point about cyclic fixpoint equations is that they may admit multiple solutions, including non-solutions to Definition 22.⁵⁹ Thus, the fixpoint equations are *underspecified*: they provide necessary, but not sufficient, conditions for a solution to equal the potentially infinite sum in Definition 22. Consider the following example.

Example 12. The following Dyna program encodes the geometric series: $[x] = 1 + a + a^2 + a^3 + \cdots$ where a = [a].

85 X += a * X.

Each term in the geometric series corresponds to the value of a derivation of x. There are infinitely many derivations for x under the program. If we plug this program into the fixpoint equations (4.2), we see that x depends on itself:

$$[\![\mathbf{x}]\!] = [\![1]\!] + [\![\mathbf{a}]\!] \cdot [\![\mathbf{x}]\!] = 1 + a \cdot [\![\mathbf{x}]\!]$$

Since we are in the real semiring, we can use subtraction and division as a shortcut to solve this equation:

$$\llbracket \mathbf{x} \rrbracket = \frac{1}{1-a}$$

Notice, however, the solution only corresponds to the sum over derivations (i.e., the value of the geometric series) when |a| < 1. The sum-over-derivations semantics for [x] when $a \ge 1$ ought to be ∞ . Reassuringly, ∞ satisfies the necessary conditions given by the fixpoint equations: $\infty = 1 + a \cdot \infty$ when a > 0. Notice,

⁵⁹In the acyclic case, however, the fixpoint equations define a unique value for each item x that equals the sum over (finitely many) derivations [x] (Definition 22).

also, that desired semantics of ∞ does not correspond to the minimal fixpoint since $\frac{1}{1-a} < \infty$. When a < -1, the series does not have a well-defined solution as it oscillates between $\pm \infty$; however, the fixpoint equations have the solution $\frac{1}{1-a}$.

Below, we give another example of underspecification that arises from having multiple solutions to a quadratic system:

Example 13. The following program encodes the Catalan series: $[\![x]\!] = c + 2c^2 + 5c^3 + 14c^4 + 42c^5 + 132c^6 + \cdots$ where $c \stackrel{\text{def}}{=} [\![c]\!]$.

Now, we consider the fixpoint equations for this program's value. We see that the fixpoint equations admit two solutions:

Assuming $0 \le c \le \frac{1}{4}$, the smaller solution $\frac{1-\sqrt{1-4}\,c}{2}$ corresponds to the sum over derivations. The larger solution $\frac{1+\sqrt{1-4}\,c}{2}$ is a fixpoint, but it is not equal to the sum over derivations.

Google Quadratic formula: $a x^2 + b x + c = 0 \land a \neq 0$ \iff $x = \frac{-b \pm \sqrt{b^2 - 4 a c}}{2a}$.

4.3 Fixpoint Iteration

Fixpoint iteration⁶¹ is a general scheme for approximating $\llbracket \cdot \rrbracket$. This scheme defines a sequence of approximate valuations $\nu^{(0)}, \nu^{(1)}, ..., \nu^{(T)}$, which approach $\llbracket \cdot \rrbracket$ as $T \to \infty$. The sequence is defined as

$$\begin{cases} \boldsymbol{\nu}^{(0)} = \bot \\ \boldsymbol{\nu}^{(t+1)} = \mathbf{T}_{\mathfrak{p}} \, \boldsymbol{\nu}^{(t)} \end{cases} \tag{4.3}$$

where \perp is the valuation for a program with no rules,⁶² and $T_{\mathfrak{p}}$ is the **step** operator for a program \mathfrak{p} :

$$(\mathbf{T}_{\mathfrak{p}} \, \boldsymbol{\nu})[\boldsymbol{x}] \stackrel{\text{def}}{=} \bigoplus_{\mathbf{R} \in \mathfrak{p}, \ (\boldsymbol{x} \oplus \boldsymbol{x} y_1 \circ \cdots \circ y_K) \in \Gamma(\mathbf{R})} \boldsymbol{\nu}[y_1] \circ \cdots \circ \boldsymbol{\nu}[y_K] \qquad \text{for } \boldsymbol{x} \in \mathbb{H}$$
 (4.4)

Equation (4.5) is simply the abstraction of the right-hand side of (4.2) from Definition 25 into a function of the valuation ν . We can rephrase the fixpoint equations (Definition 25) as a solution ν^* such that

$$\nu^* = \mathbf{T}_{\mathfrak{p}} \, \nu^* \tag{4.5}$$

Thus, like (4.2), (4.5) may admit multiple fixpoints or none at all. In practice, a suitable fixpoint ν^* is often computed by some variant of fixpoint iteration, which initializes ν to \bot , and repeatedly updates ν to $T_{\rm p} \nu$ until it converges to some $\nu^* = T_{\rm p} T_{\rm p} \cdots T_{\rm p} \bot$.

⁶¹Fixpoint iteration—in the setting of logic programming—is known as forward chaining or bottom-up evaluation (Ullman, 1988; Goodman, 1999; McAllester, 2002; Eisner, Goldlust, et al., 2005).

⁶²Note: according to Definition 22, the valuation $\{\zeta\}$ for a program with no rules $\{\zeta\}$ is not a completely empty valuation, which would map all of its arguments to the value $\underline{0}$; rather, constants and built-in relations may have non- $\underline{0}$ values under $\underline{1} = \{\zeta\}$.

⁶³For some programs, this procedure may never converge or may take infinite time to converge.

Proposition 1. The valuation estimate at time t given by fixpoint iteration $u^{(t)} = \mathbf{T}_{\mathfrak{p}}^t \perp \text{equals the total value of all derivations of height} \leq t.$

Proof (Sketch). We start with $\nu^{(0)} = \bot$. From there, we build height ≤ 1 derivations by extending the height ≤ 0 derivations (i.e., constants and builtins appearing on the right-hand sides of the rules in the program). From there, each application of the step operation combines the valuations of height $\le t$ derivations by applying the rules of the program to derive the values of all height $\le t+1$ derivations.

Fixpoint iteration is the mathematical basis for the algorithms we develop in §5. Assuming that the valuation has a convergent sum, fixpoint iteration will diligently compute the sum (up to each height) until every height has been covered; thus, it converges in the limit to the value of [.]. In that later section, we will discuss and address computational inefficiencies and provide implementation details. We will also describe an accelerated fixpoint iteration variant based on Newton's method for solving equations.

In the boolean semiring, however, there is always a unique minimal fixpoint ν^* , and for each item, x, $\nu[x]$ reaches $\nu^*[x]$ in finite time.

⁶⁴This is a form of dynamic programming: $\nu^{(t)}$ is a memo table for a version of the computation, which makes the height t an argument to break cycles.

4.4 Convergence Theory

This section briefly overviews the basic results found in the literature on semiring-weighted systems. Our goal is primarily to review; no new results are presented in this section.

An easy case where convergence is guaranteed is when the fixpoint equations are acyclic (i.e., the set of derivations is finite); here, there is no issue computing its value by fixpoint iteration or—better yet—by dynamic programming. The tricky cases are cyclic equations, as they denote sums over of infinitely many derivations, which may not be well-defined.

We have seen that *not all* fixpoints (Definition 25) are solutions to (4.1). However, fixpoint iteration (§4.3) is a procedure that sums the total value of derivations by their height (Proposition 1); thus, the fixpoint iterates approach $\llbracket \cdot \rrbracket$ assuming it is well-defined. To formalize what is meant by *approaching* a value, we require a notion of topology; we will describe a commonly used ω -continuous topology in this section.

Kuich (1997) provides sufficient conditions for when fixpoint iteration—starting from a valuation that assigns all items the value 0—converges in the limit to $[\cdot]$ (Definition 22).⁶⁵ The theory assumes that the semiring is ω -continuous (Definition 28), which fortunately includes most semirings of interest in NLP. The class of ω -continuous semirings has useful closure properties to ensure that the solution to the infinite sum over derivations equals the supre-

⁶⁵We found out about this work thanks to Goodman (1998, §2.3.1–2.3.2).

mum of the partial sums (under an appropriate ordering relation). We note that some semirings—including $\langle \mathbb{R}, +, \cdot, 0, 1 \rangle$ (unless it is restricted to $\mathbb{R}_{\geq 0}$)—are not ω -continuous—essentially because they include negative numbers, which means they can lead to undefined limits.

Before defining ω -continuous semirings, we define an ordering relation on semirings as it provides a topology on the semiring with which we can employ the usual mathematical concepts relating to convergence, such as the supremum.

Definition 26. The **natural order** (denoted by \sqsubseteq) on the members of a semiring $\langle \mathbb{V}, \oplus, \circ, \underline{0}, \underline{1} \rangle$ is defined follows: $\forall a, b \in \mathbb{V} : (a \sqsubseteq b) \iff (\exists c \in \mathbb{V} : a \oplus c = b)$. The natural order is a preorder on \mathbb{V} . However, when \sqsubseteq defines a partial order, we say that the semiring is **naturally ordered**.

Example 14 (Examples and Non-examples of Naturally Ordered Semirings).

- The natural order on real (and integer) semirings does not define a partial order: ∀a,b there always exists c = (b a) that satisfies a ⊕ c = b. Thus, ⊆ does not satisfy the antisymmetry requirement of partial order.
- The natural order on the counting semiring (or the nonnegative-real semiring) coincides with \leq : $\forall a, b$ there exists c such that $a \oplus c = b$ if and only if $b a \in \mathbb{N}$. In other words, if $a \sqsubseteq b \iff a \leq b$.
- The natural order on formal languages coincides with the \subseteq relation on sets: For $L_1, L_2 \in 2^{\Sigma^*}$, $L_1 \sqsubseteq L_2 \iff L_1 \subseteq L_2$.

Proposition 2. (Kuich, 1997, Theorem 2.1): Every naturally ordered semiring is monotonic (under \sqsubseteq), and $\underline{0}$ is its smallest element ($\forall a \in \mathbb{V} : \underline{0} \sqsubseteq a$).

Thus, for naturally ordered semirings, each semiring operation is nondecreasing on the \sqsubseteq -lattice on \mathbb{V} . The natural order can be extended to valuations ($\mathbb{V}^{\mathbb{H}}$): let $\nu, \nu' \in \mathbb{V}^{\mathbb{H}} : \nu \sqsubseteq \nu' \iff \forall x \in \mathbb{H} : \nu[x] \sqsubseteq \nu'[x]$. Next, we observe that the step operator $T_{\mathfrak{p}}$ for a given program \mathfrak{p} is monotone: $\nu \sqsubseteq \nu' \implies T_{\mathfrak{p}} \nu \sqsubseteq T_{\mathfrak{p}} \nu'$. Thus, the fixpoint iteration iterates increase \sqsubseteq on each iteration until they eventually saturate; no oscillation can cause the sum to be ill-defined. However, additional technical details are required to formally push this idea into a theorem; specifically, it is sufficient for the semiring to be ω -continuous (defined below).

Definition 27. (Kuich, 1997, pages 611–612): A semiring is **complete** if it has an infinite sum operation (\bigoplus) that satisfies the following properties for all countably infinite sets I:

(i) Base cases:

$$\bigoplus_{i\in\emptyset} a_i = 0$$
, $\bigoplus_{i\in\{j\}} a_i = a_j$, and $\bigoplus_{i\in\{j,k\}} a_i = a_j \oplus a_k$ for $j \neq k$

(ii) Infinite commutativity and associativity:

$$\bigoplus_{J\in\mathcal{J}}\bigoplus_{i\in J} a_j = \bigoplus_{i\in I} a_i$$
 where \mathcal{J} is any partitioning of I

(iii) Infinite distributivity:

$$c \circ (\bigoplus_{i \in I} a_i) = \bigoplus_{i \in I} (c \circ a_i) \ and \ (\bigoplus_{i \in I} a_i) \circ c = \bigoplus_{i \in I} (a_i \circ c) \ for \ all \ c \in \mathbb{V}$$

Definition 28. Suppose that $\langle \mathbb{V}, \oplus, \cdot, \underline{0}, \underline{1} \rangle$ is a naturally ordered semiring, and its infinite sum operator is satisfies

$$\bigoplus_{i\in\mathbb{N}} a_i = \sup\{a_1 \oplus \cdots \oplus a_i \mid i\in\mathbb{N}\} \quad \textit{for all infinite sequences } \langle a_i \rangle_{i\in\mathbb{N}} \qquad \textbf{(4.6)}$$

where the \sup^{66} is taken with respect to the ordering \sqsubseteq . If this semiring is complete, then we say that the semiring is ω -continuous.

The definitive reference for the study of these assumptions is Kuich (1997). These are the primary theoretical assumptions studied in work on semirings (e.g., Goodman (1998), Goodman (1999), Esparza et al. (2007), Green, Karvounarakis, et al. (2007), and Mohri (2009)).

Convergence-rate analysis. It is possible to study fixpoint iteration's rate of convergence when p is a contraction mapping. The limit analysis above makes weaker assumptions (i.e., ω -continuity), which are sufficient to ensure convergence in the limit, but it does not characterize the rate of convergence.

To use the contraction mapping analysis in our setting, we require that the set of valuation mappings $\mathbb{V}^{\mathbb{H}}$ is a **complete**⁶⁷ **metric space**—i.e., it must be equipped with a **distance metric** $\rho \colon \mathbb{V}^{\mathbb{H}} \times \mathbb{V}^{\mathbb{H}} \to \mathbb{R}_{\geq 0}$. For example, in the case of the real semiring, we can use the Chebyshev distance $\rho(\nu,\nu') = \max_{x \in \mathbb{H}} |\nu[x] - \nu'[x]|$. A complete metric space requires that every Cauchy sequence $\nu^{(0)}, \nu^{(1)}, \ldots$ has a limit also in $\mathbb{V}^{\mathbb{H}}$. We say $\nu^{(0)}, \nu^{(1)}, \ldots$ is a **Cauchy sequence** if, for any arbitrarily small distance r > 0, $\exists N > 0$ such that

⁶⁶The **supremum** $\sup S$ of a subset $S \subseteq P$ from a partially ordered set $\langle P, \sqsubseteq \rangle$, is the least element in P that is greater than or equal to each element of S if such an element exists. The supremum is sometimes called the **least upper bound**.

⁶⁷Note that the meaning of *complete* here is different from a *complete* semiring (Definition 27). ⁶⁸Plenty of other distance metrics are possible for the real semiring, such as the Euclidean and Manhattan distance metrics.

 $\forall m, n > N, \, \rho\left(\nu^{(n)}, \nu^{(m)}\right) < r.$ We say the step operation $\mathbf{T}_{\mathfrak{p}}$ for a given program \mathfrak{p} is a **contraction mapping** if

$$\forall \nu, \nu' \in \mathbb{V}^{\mathbb{H}} \colon \rho(\mathbf{T}_{\mathfrak{p}} \, \nu, \mathbf{T}_{\mathfrak{p}} \, \nu') \leq \lambda \cdot \rho(\nu, \nu')$$

with a **contraction rate** $\lambda < 1.^{69}$ Intuitively, a contraction mapping is a function that, when applied to any pair of points $\langle \nu, \nu' \rangle$, brings them closer together by at least a factor of λ . Note that it is *not* guaranteed that $T_{\mathfrak{p}}$ is a contraction mapping for all programs \mathfrak{p} .

If T_p is a contraction mapping, then the **Banach fixpoint theorem** guarantees that the fixpoint iteration sequence $\nu^{(0)}, \nu^{(1)}, \dots$ is Cauchy, and converges at a linear rate (as a function of n) to a unique fixpoint ν^* :

$$\rho(\nu^{(n)}, \nu^*) \le \lambda^n \cdot \rho(\nu^{(0)}, \nu^*) \tag{4.7}$$

Note that the quantity $\rho(\nu^{(0)}, \nu^*)$ in the right-hand side of (4.7) cannot be computed; fortunately, we can relax the bound so that it is computable by replacing this quantity with one that can be computed: the size of the first update $R = \rho(\nu^{(0)}, \nu^{(1)})$. This gives the so-called **prior bound**:

$$\rho(\nu^{(n)}, \nu^*) \le \frac{\lambda^n}{1 - \lambda} \cdot R \tag{4.8}$$

This bound can be rearranged to bound the number of iterations (n) that are needed to ensure that the distance to ν^* is $\leq \epsilon$:

⁶⁹Contraction mappings are a special case of Lipschitz-continuous functions where the function maps onto itself (i.e., the functions domain and range the same). Thus, the constant λ is sometimes called a Lipschitz constant.

$$n \ge \frac{\log\left(\frac{\epsilon(1-\lambda)}{R}\right)}{\log \lambda} \implies \rho(\nu^{(n)}, \nu^*) \le \epsilon$$
 (4.9)

Another useful bound is the **posterior bound**:

$$\rho(\nu^{(n)}, \nu^*) \le \frac{\lambda}{1 - \lambda} \cdot \rho(\nu^{(n)}, \nu^{(n-1)}) \tag{4.10}$$

The posterior bound relates the update's size at iteration n to the remaining distance to ν^* . The posterior bound shows that if the most recent update (at iteration n) is small, we are close to the solution. Thus, the posterior bound gives an approximate termination condition.

$$\frac{\lambda}{1-\lambda} \cdot \rho(\nu^{(n)}, \nu^{(n-1)}) \le \epsilon \quad \Longrightarrow \quad \rho(\nu^{(n)}, \nu^*) \le \epsilon \tag{4.11}$$

In practice, however, we do not know λ , so the threshold is just on $\rho(\nu^{(n)}, \nu^{(n-1)})$.

Summary

This section formalized the semantics of Dyna programs as a sum over semiring-weighted derivations. We described fixpoint iteration, the mathematical basis of the execution algorithms we will develop in § 5. We provided an overview of the sufficient conditions for which fixpoint iteration converges to [].

Chapter 5

Inference Algorithms

This section develops general algorithms for performing inference in Dyna programs (i.e., evaluating $[\cdot]$). These general algorithms are why Dyna programs do not require the specification of implementation details: programmers provide a specification of the logic, and the Dyna system takes care of the control flow. We describe several algorithms:

- In §5.1, we describe three implementations of fixpoint iteration of increasing sophistication: naïve (§5.1.1), semi-naïve (§5.1.2), and prioritized (§5.1.3).
- In §5.2, we then describe an *accelerated* fixpoint algorithm, which is based on a version of Newton's method for solving fixpoint equations over a semiring (Esparza et al., 2007).
- In §5.3, we describe a specialized solver for certain linearly recursive programs. This specialized solver is leveraged as a subroutine by the accelerated algorithm.

• In §5.4, we describe a unified, hybrid algorithm that solves large programs by breaking them up into smaller subprograms that can be solved separately with the best applicable method.

Solution representations. Recall from §2 that we represent the answer to a query as a program of a restricted form where each of its rules is simple. A rule is simple if its body contains only constant values and builtins. We call a program consisting of only simple rules a **simple program**. In this chapter, we use simple programs to represent the valuation function $[\cdot]$. This allows us to compactly represent the solutions to certain programs with infinite valuations, such as⁷⁰

$$\begin{cases} 88 & x(I) += (I > 0). \\ 89 & x(I) += 0.5 * x(I). \end{cases} \xrightarrow{\text{solve}} \begin{cases} 90 & x(I) += 2 * (I > 0). \end{cases}$$

This generalizes prior work by Eisner, Goldlust, et al. (2005), which required solutions to be represented as a hash table mapping ground terms to their values. In our notation, that translates to a restriction that each rule in the solution has a ground head, and its body is a constant value. The consequence of their assumption is that they can only represent the solutions to programs with finitely many non-0 items.

We will use simple programs in our algorithms to represent intermediate valuation functions. We will use the result-stream manipulation operations (developed in §3.1) in the algorithms of this section, as they make working it possible to manipulate these program representations seamlessly.

⁷⁰We provide additional examples in §A.

5.1 Algorithms for Fixpoint Iteration

In §4.3, we discussed the mathematical foundation for fixpoint iteration. Here, we will describe the algorithm procedurally and discuss implementation details. In the logic programming literature, fixpoint iteration is known under the names bottom-up evaluation or forward chaining. We will give three algorithms of increasing sophistication: naïve (§5.1.1), semi-naïve (§5.1.2), and prioritized (§5.1.3). We do not recommend using the algorithms naïve (§5.1.1) and semi-naïve (§5.1.2); we have included them for pedagogical purposes to motivate the prioritized algorithm (§5.1.3). Our version of prioritized fixpoint iteration generalizes the forward chaining algorithm of Eisner, Goldlust, et al. (2005) to support non-range-restricted rules, built-in relations, and delayed constraints.

5.1.1 Naïve Fixpoint Iteration

This section describes our simplest implementation of fixpoint iteration. In the pseudocode provided in Algorithm 6, we can see that the algorithm is a remarkably straightforward translation of the mathematical definition of (4.4) thanks to our abstractions for manipulating streams (§3.1).⁷³

In the code, we initialize the **chart** m to the empty program \emptyset (i.e., the program with no rules) so its valuation equals \bot , as in (4.3). We then repeatedly

⁷¹E.g., Ullman (1988), Goodman (1999), McAllester (2002), and Eisner, Goldlust, et al. (2005).

⁷²Our version is most like Eisner, Goldlust, et al. (2005), as it has been generalized to support semiring values.

 $^{^{73}}$ These abstractions allow us to specify each of our inference algorithms *as if* the program's rules do not contain variables.

Algorithm 6 naive: Solve program p by (naïve) fixpoint iteration

```
1. def naive(p):
        ▷ This procedure implements fixpoint iteration (4.3) on the function stepp.
        \mathfrak{m} \leftarrow 75
                                                \triangleright \{ \} is a program with no rules, its valuation equals \perp as in (4.3).
        while true:
             \mathfrak{m}' \leftarrow \mathsf{step}_{\mathfrak{p}}(\mathfrak{m})
                                                                                                                        \triangleright Apply the step operation
            if \mathfrak{m} = \mathfrak{m}': return \mathfrak{m}
             \mathfrak{m} \leftarrow \mathfrak{m}'
s. \triangleright The procedure \mathsf{step}_\mathfrak{p}(\mathfrak{m}) implements T_\mathfrak{p}(\nu) from (4.4); here, the valuation \nu is represented as
      a program m.
9. def step<sub>n</sub>(\mathfrak{m}):
        \mathfrak{m}' \leftarrow 75
        for (x \oplus = y_1 \circ \cdots \circ y_K) \in \mathfrak{p}:
            \mathbf{m}'[\mathbf{z}] \oplus = \mathbf{m}[\mathbf{y}_1] \circ \cdots \circ \mathbf{m}[\mathbf{y}_K]
        return m'
```

apply the step function $step_p$, which implements T_p from (4.4), until a fixpoint is reached (i.e., the condition on line 6.6 is satisfied).

We see that $\operatorname{step}_{\mathfrak{p}}(\mathfrak{m})$ implements $\mathbf{T}_{\mathfrak{p}}(\mathfrak{m}[\![\cdot]\!])$ as it sums the contributions of each rule in the program $(x \oplus = y_1 \circ \cdots \circ y_K) \in \mathfrak{p}$ by instantiating its body against the current chart \mathfrak{m} with the help of the result-stream abstractions: $\mathfrak{m}[y_1] \circ \cdots \circ \mathfrak{m}[y_K]$. Line 6.12 adds the contributions of the current rule to the next iterate \mathfrak{m}' using $\mathfrak{m}'[x] \oplus = \mathfrak{m}[y_1] \circ \cdots \circ \mathfrak{m}[y_K]$. It is straightforward to verify that \mathfrak{m} and \mathfrak{m}' are always simple programs.

We note that the runtime of Algorithm 6 is not as good as it could be. In the following sections, we will improve upon it in several ways.

Termination condition. Our algorithm's termination condition (line 6.6) is $\mathfrak{m} = \mathfrak{m}'$. Termination can be guaranteed for the following cases:

• Finite acyclic programs (regardless of the semiring).

• Finite programs over a semiring that is guaranteed to saturate after a finite number of iterations. The family of *K*-closed semirings (Definition 19) is a setting where this is always true; this family includes many important semirings such as the boolean semiring and the optimization semirings (§3.3). However, it does not include many important semirings, such as the real semiring.

Approximate termination condition. Defining an approximate termination condition to replace line 6.6 depends on the specific semiring.⁷⁴ A natural choice for the real semiring is $\forall x \colon |\mathfrak{m}[x] - \mathfrak{m}'[x]| \le \epsilon$ for some small threshold $\epsilon \ge 0$. For other semirings, |v-v'| should be replaced with an appropriate distance metric $\rho(v,v')$. This choice is inspired by the posterior bound (4.10) in the Banach fixpoint theorem (discussed in §4.4). The number of fixpoint iterations required to reach the threshold ϵ may be difficult to bound unless the step function is a contraction mapping (§4.4).

5.1.2 Semi-Naïve Fixpoint Iteration

The seminaive algorithm is a more efficient implementation of Algorithm 6. The improvement is based on the observation that on iteration h+1, many of the values for items in $\mathfrak{m}^{(h)}$ are unchanged; however, the step operator *rederives* them for $\mathfrak{m}^{(h+1)}$. The seminaive algorithm refines the induction on height by precisely deriving the set of derivations with height =h rather than $\leq h$ on

74Our implementation has hooks for defining approximate equality for a given semiring.

each step. Doing so improves efficiency, as each step does not need to rederive derivations of lower heights.

Pseudocode is given in Algorithm 7. The algorithm's structure is similar to Algorithm 6, except that in addition to \mathfrak{m} (and \mathfrak{m}'), we maintain the values of new heights \mathfrak{d} (and \mathfrak{d}'). It is straightforward to verify the invariant that \mathfrak{m} , \mathfrak{d} , \mathfrak{m}' , and \mathfrak{d}' are simple programs. We argue for the correctness of the algorithm in Proposition 3.

Algorithm 7 seminaive: Solve program p by semi-naïve fixpoint iteration

```
1. def seminaive(\mathfrak{p}):
           \mathfrak{m} \leftarrow 75
           \mathfrak{d} \leftarrow \mathsf{step}_{\mathfrak{p}}(\mathfrak{m})
                                                                                                                                        ▷ Initialize with step from Algorithm 6
           while true:
                  \langle \mathfrak{m}', \mathfrak{d}' \rangle \leftarrow \mathsf{step}'_{\mathfrak{p}}(\mathfrak{m}, \mathfrak{d})
                  if \mathfrak{m} = \mathfrak{m}': return \mathfrak{m}
                  \langle \mathfrak{m}, \mathfrak{d} \rangle \leftarrow \langle \mathfrak{m}', \mathfrak{d}' \rangle
8. \triangleright Incremental step function; applies change \eth to chart \mathfrak{m} and derives a new change \eth'.
9. def step'<sub>n</sub>(\mathfrak{m}, \mathfrak{d}):
                                                                        Should omit constants and builtins;
           ð' ← ? S <
                                                                       these are added by line 7.3.
           \mathfrak{m}' \leftarrow \mathfrak{m} \oplus \mathfrak{d}
                                                                                                                                                                                   \triangleright Accumulate changes
          \begin{array}{l} \textbf{for}\; (x \oplus = y_1 \circ \cdots \circ y_K) \in \mathfrak{p} \colon \\ \textbf{for}\; k \in [1{:}K]{:}^{75} \end{array}
                        {\color{red}\mathfrak{d}'[\mathbf{\textit{x}}] \oplus = \mathfrak{m}'[\mathbf{\textit{y}}_{[:k)}] \circ \textcolor{red}\mathfrak{d}[\mathbf{\textit{y}}_{k}] \circ \mathfrak{m}[\mathbf{\textit{y}}_{(k::]}]}
           return \langle \mathfrak{m}', \mathfrak{d}' \rangle
```

Proposition 3. On iteration h of Algorithm 7 (line 7.7), the value of \mathfrak{m} is the valuation function for heights < h, and the value of \mathfrak{d} is the valuation function restricted height = h.

Proof (Sketch). According to Definition 24, we obtain a=h+1 derivation when we apply one rule to a collection of at least one =h derivation and any number of $\overline{^{75}}$ For efficiency, the join algorithm should use $\mathfrak{d}[y_k]$ as its outermost loop.

< h derivations. To illustrate, consider the case of a derivation with 2 children. A rather direct way to partition the events that create height = h+1 derivations is the following three cases, which identify all positions where the = h subderivation can appear:

$$\operatorname{height}\left(\begin{array}{c} x\\ \\ y_1 \\ \\ \end{array}\right) = \operatorname{height}\left(\begin{array}{c} x\\ \\ y_1 \\ \\ \end{array}\right) = \operatorname{height}\left(\begin{array}{c} x\\ \\ y_1 \\ \\ \end{array}\right) = h+1$$

However, a more efficient—albeit less obvious—partitioning of the event above merges the first two cases into one case ($\leq h$)—giving two cases instead of three:

$$\operatorname{height}\left(\begin{array}{c} x\\ y_1 & y_2\\ \\ & = h \end{array}\right) = \operatorname{height}\left(\begin{array}{c} x\\ y_1 & y_2\\ \\ & = h \end{array}\right) = h+1$$

What we have done here is to identify events with their rightmost =h subderivation.⁷⁶ We depict the general case below (for k=1...K):

We refer to this strategy for grouping events as the **cursor trick**. We refer to the position of the rightmost =h subtree as the *cursor* as it moves from left to right in a given rule body to control which version of each subgoal's value $(\mathfrak{m}, \mathfrak{d}, \mathfrak{m}')$ is used in the update.

⁷⁶When we consider derivations with more than two children, the benefit of the latter decomposition becomes more apparent. In the latter decomposition, there are K cases to consider, whereas the former has $\mathcal{O}(2^K)$ because it considers each subset of positions where the maximal height can appear.

Algorithm 7 operationalizes this construction as follows. The function step' derives the values for the set of derivations of the next height (=h+1) by applying each rule in the program once: for a given rule $x \oplus = y_1 \circ \cdots \circ y_K \in \mathfrak{p}$, we consider plugging in each item from \mathfrak{d} as the rightmost =h child. We do so by looping over values of the cursor k from 1 to K (line 7.13), and applying the update on line 7.14 (repeated below):

$$\mathfrak{d}'[x] \oplus = \mathfrak{m}'[y_1] \circ \cdots \circ \mathfrak{m}'[y_{k-1}] \circ \mathfrak{d}[y_k] \circ \mathfrak{m}[y_{k+1}] \circ \cdots \circ \mathfrak{m}[y_K]$$

We see that the update has the same pattern as (5.1), which is why \mathfrak{d}' corresponds to values contributed by the set of the derivation of height = h+1.

Example 15. We give a simple example of the cursor trick in action.

We begin by initializing $\mathfrak{m}=\bigcap and \mathfrak{d}= step_{\mathfrak{p}}(\mathfrak{m}).$ Then, on the first iteration:

$$\mathfrak{m}=\text{isom} \ \mathfrak{d}=\begin{cases} y \text{ += 2.} \\ z \text{ += 1.} \end{cases} \quad \overset{\text{step'}}{\Longrightarrow} \quad \mathfrak{m'}=\begin{cases} y \text{ += 2.} \\ z \text{ += 1.} \end{cases} \quad \mathfrak{d'}=\begin{cases} x \text{ += 2 * 2 * 2 * 2.} \\ y \text{ += 1.} \end{cases}$$

Next, we see the cursor trick in action in \mathfrak{d}' :

$$\mathfrak{m} = \begin{cases} y \ += \ 2. \\ z \ += \ 1. \end{cases} \qquad \mathfrak{d} = \begin{cases} x \ += \ 8. \\ y \ += \ 1. \end{cases} \qquad \overset{step'}{\Longrightarrow} \qquad \mathfrak{m}' = \begin{cases} x \ += \ 8. \\ y \ += \ 3. \\ z \ += \ 1. \end{cases} \qquad \mathfrak{d}' = \begin{cases} x \ += \ 1 \ * \ 2 \ * \ 2. \\ x \ += \ 3 \ * \ 1 \ * \ 2. \\ x \ += \ 3 \ * \ 3 \ * \ 1. \end{cases}$$

Next, we apply the update and see that x has the correct value of 27 in m':

$$\mathfrak{m} = \begin{cases} x += 8. \\ y += 3. \\ z += 1. \end{cases} \quad \mathfrak{d} = \begin{cases} x += 19. \end{cases} \quad \overset{\text{step'}}{\Longrightarrow} \quad \mathfrak{m'} = \begin{cases} x += 27. \\ y += 3. \\ z += 1. \end{cases} \quad \mathfrak{d'} = \emptyset$$

The seminaive *halts* on the next iteration with the final chart:

$$m = \begin{cases} x & += 27. \\ y & += 3. \\ z & += 1. \end{cases}$$

Efficiency. In logic programming, values are always boolean; thus, once a *single* derivation for a given item is found, we no longer need to consider expanding. In other words, once an item has a single derivation, its value has *converged*—i.e., it has reached its fixpoint. Thus, the seminaive strategy is particularly efficient in the boolean semiring because it only does work proportional to the change m, and it does not perform any unnecessary rule application. However, when considering general semiring values, we see that seminaive algorithms can be inefficient. In particular, seminaive will combine values for items before they have converged, and it will repeat the combination each time a taller derivation of an item is discovered.

Example 16. Below is an example family of acyclic programs with a linear size in a parameter n. On this family of programs, semi-naive runs in $O(n^2)$ time due to unnecessary repropagation.⁷⁸

```
121 x(I+1) += x(I), I < n.

122 x(1) += y(I).

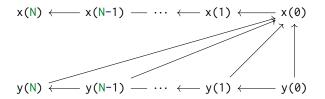
123 y(I+1) += y(I), I < n.

124 y(1) += 1.
```

We can visualize the direct dependencies among the items as a graph

⁷⁷This requires that the loop over rules in line 7.12 is modified so that only rules with a non-<u>0</u> grounding are enumerated.

⁷⁸We are using a short-hand that allows certain built-in relations, e.g., I+I, to be written in place: x(I+1) += x(I), I < n translates to x(I') += x(I), I' is I+1, I < n.



At a high level, the reason why seminaive runs in quadratic time is because $x(\emptyset)$'s value will update $\mathcal{O}(n)$ times and, thus, items that depend on $x(\emptyset)$ will update $\mathcal{O}(n)$ times as well. After a warm-up period, \emptyset will have $\mathcal{O}(n)$ values propagated per iteration. These big changes happen $\mathcal{O}(n)$ times before converging, which gives the quadratic time.

In the next section, we will develop another algorithm that allows us to compute the value of each item in this program in topologically sorted order—achieving a much-improved runtime of $\mathcal{O}(n)$ on this example.

5.1.3 Prioritized Fixpoint Iteration

In this section, we will give a more flexible strategy for semi-naive evaluation, which uses an agenda (work list) to safely eliminate height as a synchronization barrier. The benefit is that we can avoid *repropagation*: rather than advancing all items from height $\leq h$ to height $\leq h+1$ in lockstep, we can allow some items to advance way beyond the others. Ideally, we would defer propagation until each item reaches its fixpoint value. Once a set of items converges, we do not need to continue combining with other converged items (we only need to combine a set of converged items once). In cyclic graphs, blocking updates until a fixpoint is reached is impossible (deadlock), but we can focus updates on propagating

Algorithm 8 agenda: Solve program p by agenda-based fixpoint iteration

```
1. def agenda(p):
          \mathfrak{m} \leftarrow ? S; \mathfrak{m}' \leftarrow ? S
          \mathfrak{d} \leftarrow \mathsf{step}_{\mathfrak{p}}(\mathfrak{m})
                                                                                                             ▷ Initialize with step from Algorithm 6
          while \emptyset \neq \emptyset: \frown May contain delayed constraints
               (y' \oplus = \overset{\smile}{v}) \leftarrow 0.\mathsf{pop}()
                                                                                                                   Pick any update from ∂ to apply
               \mathfrak{m}'[y'] \oplus = v
               if \mathfrak{m}'[y'] = \mathfrak{m}[y']: continue

▷ Update has no effect; no need to propagate<sup>80,81</sup>

               for (x \oplus \exists y_1 \circ \cdots \circ y_K) \in \mathfrak{p}:
                    for k \in [1:K]:
                         \begin{array}{c} \mathbf{0}[\mathbf{z}] \oplus \mathbf{m}'[\mathbf{y}_{::k})] \circ v[\mathbf{y}_k = \mathbf{y}'] \circ \mathbf{m}[\mathbf{y}_{(k:]}] \end{array}
10
               \mathfrak{m}[y'] \oplus = \underline{v}
                                                                                                                                                           \triangleright Now, \mathfrak{m} = \mathfrak{m}'
11
          return m
12
```

within a cyclic subgraph before moving on. In other words, we can prioritize updates within a strongly connected component (Definition 12), allowing it to reach a fixpoint before advancing to the next strongly connected component.⁷⁹

The key to allowing this generalization is to generalize of from Algorithm 7 to buffer changes before they are committed to m. This allows some items to advance their valuations to larger heights while others are deferred. The values accumulated in of correspond to the sums over heights that have yet to be accumulated into m—it is no longer limited to a single next height. Assuming that no updates are deferred indefinitely, the agenda strategy's valuation estimate should converge to [.] (under the same conditions that seminaive would converge).

The policy for determining which **update** to pop from the **agenda** ∂ (line 8.5) is called the **prioritization policy**. It is typically implemented using a priority queue. We provide guidance on prioritization function design later in §5.1.4.

⁷⁹Goodman (1998) calls SCCs "looping buckets."

Pseudocode for our algorithm is provided in Algorithm 8.82 We note that Algorithm 8 is based closely on the agenda-based algorithm given in Eisner, Goldlust, et al. (2005). We have made subtle but necessary modifications to their algorithm to support nonground updates.⁸³ Much like Algorithm 7, the agendabased algorithm uses the cursor trick, 84 which—at least conceptually—requires multiple versions of the value of the **driver** item y' that is popped from the agenda (\mathfrak{d}): the new version $\mathfrak{m}'[y']$, the change $\mathfrak{d}[y']$, and the old version $\mathfrak{m}[y']$. We call the queries $\mathfrak{m}'[y_{[:k)}]$ and $\mathfrak{m}[y_{(k:]}]$ on line 8.10 **passenger** queries. If the y' does not match multiple subgoals, then the passenger queries will see the values with the most up-to-date value that is not pending on the change buffer. The update $y' \oplus v$ is only *committed* to m on line 8.11 after the update has been carefully pushed through the rules in the program (using the cursor trick to avoid double counting the update); at this point $\mathfrak{m} = \mathfrak{m}'$. The strategy used by Eisner, Goldlust, et al. (2005) is essentially the same. However, their strategy is only correct for ground updates: it does not extend to the case of nonground

⁸⁰Note that it is safe to allow false negatives on this equality test on line 8.7, but not false positives. However, false negatives may result in more iterations for 0 to become empty (line 8.4).

⁸¹The equality test on line 8.7 is between the results of the query specified by y' in each of \mathfrak{m}' and \mathfrak{m} . When y' is nonground, this test may no longer run in constant time as we are comparing the equality of result streams.

⁸²Much like Algorithm 7 is straightforward to verify the invariant that m, d, m', and d' are simple programs.

⁸³The key difference is that we represent the chart and agenda as programs, whereas Eisner, Goldlust, et al. (2005) require them to each be a hash table mapping ground items to semiring values. This representation difference allows us to represent certain infinite valuation functions efficiently. A more superficial difference is that we maintain two versions of the chart m and m', whereas Eisner, Goldlust, et al. (2005) represent the old value of the driver (i.e., that in m) using a local variable. A nonground update would cause Eisner, Goldlust, et al. (2005)'s pseudocode to loop over the infinite Herbrand universe matching the update; their implementation, on the other hand, would throw an exception in this case.

⁸⁴A version of the cursor trick appears in Eisner, Goldlust, et al. (2005).

5.1.4 Prioritization Functions

We use the prioritization policies in each iteration of Algorithm 8 (line 8.5) to determine which update to propagate next. In principle, any prioritization policy that does not *starve* any item from popping results in convergence (in the limit) to a fixpoint—assuming that fixpoint iteration converges at all. A simple default policy with this property is first-in, first-out (FIFO); however, FIFO is essentially just the seminaive algorithm, which we saw in Example 16 can be unnecessarily slow to converge.

This section describes a particular form of prioritization policy called a prioritization function. The way it works is that the prioritization function assigns each update in \mathfrak{d} a numeric priority at the time the update is *pushed* (i.e., when it is added to \mathfrak{d} on either line 8.3 or line 8.10).⁸⁶ When we call $\mathfrak{d}.pop()$ on line 8.5, the highest-priority update (with ties broken arbitrarily) is returned, and it is removed from \mathfrak{d} . The priority values of each update can be stored in a priority queue data structure that indexes the updates in \mathfrak{d} to make these push and pop operations efficient.

An approach to prioritization proposed by Goodman (1998) first runs the program in the boolean semiring to estimate the set of items $\overline{\mathfrak{m}}$ that have non-0

⁸⁵ Please refer to §A for examples involving nonground updates and delayed constraints.

⁸⁶It is also possible to assign items non-numeric scores; the only requirement is that the priority values are partially ordered.

Algorithm 9 support: Estimate an upper bound on the support set of p.

```
1. def support(p):
        \overline{\mathfrak{p}} \leftarrow \mathsf{booleanize}(\mathfrak{p})
                                                        \triangleright Convert \mathfrak{p} into a program \overline{\mathfrak{p}} in the boolean semiring by
                                                           replacing each non-0 constant in p with true. 88
        \overline{\mathfrak{m}} \leftarrow \operatorname{solve}(\overline{\mathfrak{p}})
                                                  \triangleright Warning: May not terminate when \overline{\mathfrak{p}} has infinitely many
                                                      true items (see text).
       return disjoin(\overline{\mathfrak{m}})
                                                                                 \triangleright Make sure \overline{\mathfrak{m}}'s rules denote disjoint sets.
5. def disjoin(\overline{\mathfrak{m}}):
       \triangleright Transform \overline{\mathfrak{m}} into a program whose rules denote nonoverlapping sets
       \triangleright This procedure can be sped up by partitioning \overline{\mathfrak{m}} into hash buckets where items
           in the same bucket might overlap, and those in different buckets cannot.
        while true:
            find r, s \in \overline{\mathbb{m}} such that r \cap s and r \neq s; if no such pair exists, break

ightharpoonup Replace the overlapping items with their tightest upper bound (merge). ^{89,90}
           \overline{\mathfrak{m}} \leftarrow (\overline{\mathfrak{m}} \setminus \{r, s\}) \cup \{\mathsf{merge}(r, s)\}
       return m
12.
```

value.⁸⁷ Algorithm 9 provides pseudocode this approach; it has been generalized beyond Goodman (1998)'s setting to allow nonground items.

Here any solver (solve) can be used to solve \overline{p} , such as seminaive, or agenda (with any nonstarving prioritization function). However, for this method to be effective, we need solve to terminate.⁹¹ Fortunately, we only require an *upper bound* on the support to move forward. In §6.1, we describe a "relaxed" forward chaining method that can be used to estimate a upper bound \overline{m} in finite time.

⁸⁷In 6.1, we further discuss the use of "booleanization" to estimate the set of non-0 items.

⁸⁸Here booleanize guarantees that $\forall x \in \mathbb{H} : \mathfrak{p}[x] \neq 0 \Longrightarrow \overline{\mathfrak{p}}[x] = \text{true}.$

⁸⁹ Here merge is the anti-unification algorithm (Reynolds, 1970; Plotkin, 1970). For example, the pair x = f(3, Y), y = f(X, 4) is rounded up to z = f(X, Y).

⁹⁰Using merge loses precision as it finds an *upper bound* on the original $\overline{\mathbb{m}}$ that is disjoint. In principle, we could completely remove the imprecision introduced by merge. By replacing line 9.11 with $\overline{\mathbb{m}} \leftarrow (\overline{\mathbb{m}} \setminus \{r, s\}) \cup \{r \cap s, r \setminus s, s \setminus r\}$. However, this choice has the computational challenge of efficiently representing differences of sets of nonground terms (Filardo, 2017, chapter 4); we use merge in practice.

⁹¹In some cases, a booleanization may drastically overestimate the set of items that forward chaining needs to care about. For example, a program like goal += p(N). p(z). p(s(N)) += 0.5 * p(N). has infinitely many nonzero items, but the value of $p(s^k(z))$ is very small as a function of k: $[p(s^k(z))] = 0.5^k$ so the values go to 0 quickly. So, we obtain a good estimate of goal with just a few iterations of fixpoint iteration. However, booleanization enumerates an infinite number of items since $[p(s^k(N))] = true$ for all $k \in \mathbb{N}$.

Algorithm 10 incoming and outgoing: Lazily enumerate incoming and outgoing arcs in program p's dependency graph given its estimated support $\overline{\mathfrak{m}}$.

```
 \begin{array}{lll} & \verb| Lazily enumerate head-to-subgoal arcs \\ & \textbf{def} \ \mathsf{incoming}(\pmb{x}) \colon & \textbf{def} \ \mathsf{outgoing}(\pmb{y}') \colon \\ & \textbf{for} \ (y_1 \circ \cdots \circ y_K) \in \mathfrak{p}[\pmb{x}] \colon & \textbf{for} \ (\pmb{x} \oplus = y_1 \circ \cdots \circ y_K) \in \mathfrak{p} \colon \\ & \textbf{for} \ _{-} \in \overline{\mathfrak{m}}[y_1 \circ \cdots \circ y_K] \colon & \textbf{for} \ k \in [1:K] \colon \\ & \textbf{for} \ k \in [1:K] \colon & \textbf{for} \ _{-} \in \overline{\mathfrak{m}}[y_{[:k)} \circ (y_k = y') \circ y_{(k:]}] \colon \\ & \textbf{yield} \ x & \textbf{yield} \ x \\ \end{array}
```

From the estimated set of non- $\underline{0}$ -valued items $\overline{\mathfrak{m}}$, we can find a topologically ordered list of strongly connected components of the item–item dependency graph, $G = \langle V, E \rangle$ where V is the set of items in $\overline{\mathfrak{m}}$, and E is the set of edges given by the head-to-subgoal relationship in the instantiated program. Algorithm 10 describes how to enumerate the incoming (and outgoing) arcs lazily. 92,93

The prioritization function is then the mapping from each item to the (integer) position of its SCC in the reverse toposort. Determining the toposorted SCCs with Algorithm 5 runs in linear time in the instantiated program's size $\mathcal{O}(|V| + |E|)$ using linear space in the number of items $\mathcal{O}(V)$. This same SCC analysis will also be used in our hybrid algorithm (§5.4).

⁹²For completeness, we have also included outgoing.

⁹³Note that the pseudocode implicitly relaxes any delayed constraints coming from $\overline{\mathfrak{m}}$ because it ignores them in the "**for** $\underline{} \in \overline{\mathfrak{m}}[...]$ " loops. This may overestimate the set of arcs in G, and consequently, the SCC analysis will be coarse. We also assume that builtins *delay* rather than throw instantiation faults (§3.1.2.2).

⁹⁴Goodman (1998) missed a trick when they first proposed this scheme: the graph edges do not need to be materialized to run Algorithm 5. Specifically, we can use incoming (or outgoing) to enumerate the dependents of a node lazily. This procedure only requires constant space because the queries and unification operations performed therein only require constant space using our result-stream implementation (§3.1.2.1). Simultaneously with this dissertation, Eisner (2023) published a detailed overview of algorithms for efficient inference in weighted deduction systems. Both our approach and theirs leverage the observation above to keep space complexity linear in |V|. Their "weighted cyclic forward-chaining algorithm" is closely related to our agenda with SCC-based prioritization. The difference is that our agenda algorithm coordinates all of its control flow through an agenda; they use a recursive algorithm for recursive solving of each SCCs (like our hybrid algorithm (§5.4)). Their approach considers weighted deduction systems that are not limited to semiring weights.

We note that if the values of only a subset of items targets $\subseteq \mathbb{H}$ are required, then we may wish to use the toposorted SCCs of sccs(incoming, targets), as it will only enumerate SCCs that targets depend on. Skipping SCCs that are not needed by the targets nodes can provide large savings in cases where the omitted SCCs are slow to converge, or the semiring operations are expensive.

To avoid additional factors in the runtime, we can use a bucket queue (Dial, 1969) to implement the pop operation in this case efficiently. With this prioritization scheme, we can solve Example 16 using Algorithm 8 in linear time, which is considerably better than the quadratic runtime of Algorithm 7. Notice that if G contains cycles, each cycle is "solved" before moving to the next. If convergence is not guaranteed (e.g., in the real semiring), this can create a barrier to progress on other program components (i.e., starvation of updates).

Discussion. Below are some alternative techniques for specifying a prioritization function:

- User-customizable prioritization functions (Eisner and Filardo, 2011).
- Synthesize a prioritization function using static analysis and satisfiability modulo theories (Colón and Sipma, 2001; Podelski and Rybalchenko, 2004; Bagnara et al., 2010).
- Lastly, we mention that there may be benefits to using carefully designed prioritization *within* an SCC. In other related settings, prioritization has been shown to improve convergence rates.⁹⁵ Future work may want to investigate

⁹⁵ For example, solvers for Markov decision processes (see Wingate and Seppi (2005) for a survey),

methods based on minimum feedback arc sets, which may be regarded as an approximate topological ordering when cycles are present (e.g., Baharev et al. (2016)), and machine learning of prioritization heuristics (e.g., Jiang et al. (2012) and Vieira, Francis-Landau, et al. (2017)).

5.2 Newton's Method

Our newton method is an accelerated solver for fixpoint problems, based on Esparza et al. (2007)'s generalization of Newton's algorithm for solving fixpoint equations over a semiring. What is interesting about newton relative to the fixpoint iterations methods is that it can sum infinitely many derivations in a single iteration. More specifically, each iteration sums the set of derivations by increasing *dimension* (defined below) rather than by its height as in fixpoint iteration (Proposition 1).⁹⁶

Etessami and Yannakakis (2009) showed that fixpoint iteration could converge exponentially slower than Newton's method.

loopy belief propagation (see Elidan et al. (2006) for a representative method), nonlinear systems of equations (see Baharev et al. (2016) for a survey), and for inference in superior semirings (e.g., Felzenszwalb and McAllester (2007)).

⁹⁶In our setting, we use newton to sum over semiring-weighted derivations; thus, our presentation will differ from that of Esparza et al. (2007)'s which takes Newton's method as their starting point. Newton's method for solving a fixpoint equation $\mathbf{x} = \mathbf{f}(\mathbf{x})$ over the reals where $\mathbf{x} \in \mathbb{R}^D$ and $\mathbf{f} : \mathbb{R}^D \to \mathbb{R}^D$ is an iterative scheme with the update rule: $\mathbf{x}^{(t+1)} \leftarrow \mathbf{x}^{(t)} + (\mathbf{I} - \mathbf{J}^{(t)})^{-1}(\mathbf{f}(\mathbf{x}^{(t)}) - \mathbf{x}^{(t)})$ where $\mathbf{J}^{(t)}$ is the Jacobian matrix of partial derivatives of \mathbf{f} evaluated at $\mathbf{x}^{(t)}$. Generalizing to the semiring setting requires retrofitting Newton's method to accommodate semirings that do not have subtraction, division, or differentiation. For our purposes, such retrofitting would be a distraction; we refer the reader to any of the following works, which discuss the connection: Esparza et al. (2007), Reps et al. (2017), and Luttenberger (2009). We found the presentation in Schlund (2016) the clearest.

Example 17. Consider the following program based on Etessami and Yannakakis (2009, Theorem 3.2). This program has a single fixpoint [x] = 1.97

Etessami and Yannakakis (2009) show that fixpoint iteration is exponentially slower than Newton's method in this example. Specifically, fixpoint iteration requires roughly 2^k iterations to gain k bits of precision, whereas Newton's method gains one bit of precision per iteration.

Definition 29. The **dimension** $dim(\delta)$ of a derivation δ is defined inductively as

- (i) Base case: $\dim(\delta) \stackrel{\text{def}}{=} 0$ if δ is a value or builtin
- (ii) If δ is of the form $x \oplus \delta_1 \circ \cdots \circ \delta_K$ where no δ_k is a product
 - If K = 0: $dim(\delta) = 0$.
 - If $K \ge 1$: Let $d = \max(\dim(\delta_1), \dots, \dim(\delta_K))$. If the argmax is unique (i.e., no ties for the max), $\dim(\delta) \stackrel{\text{def}}{=} d$. Otherwise, $\dim(\delta) \stackrel{\text{def}}{=} d + 1$.
- (iii) If δ a product of the form $\delta_1 \circ \cdots \circ \delta_K$: $\dim(\delta)$ is not defined for products because the tie-checking definition is not associative.

Please refer to Figure 5.1 for an illustration.

The dimension of a derivation is always \leq its height, but it is often \ll . The set of derivations with dimension d can be infinite. In contrast, the set of derivations with fixed height is finite. 98 For example, each derivation in the geometric series

⁹⁸The statement about fixed height derivations assumes ignores the possibility that a built-in relation creates an infinite number of derivations with height =h+1.

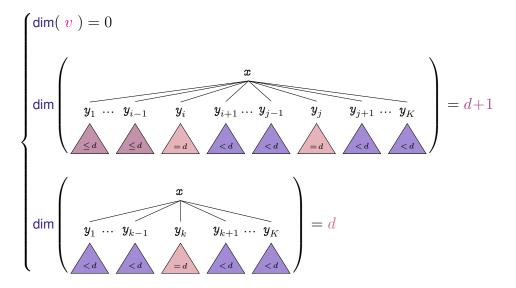


Figure 5.1: Illustration of the dimension of a derivation (Definition 29).

program (Example 10) has dimension ≤ 1 . What is remarkable about the set of derivations with dimension d+1 is that there is an efficient *linear* structure relating it to the sets of derivations with dimensions $\leq d$, = d, and < d. This allows us to grow the set of derivations we have summed over much more quickly using a specialized linear solver (§5.3).

We provide pseudocode for newton in Algorithm 11.⁹⁹ The high-level picture of newton is much like seminaive forward chaining: newton tracks its analogs of \mathfrak{m} and \mathfrak{d} , except here \mathfrak{m} is the set of dimension < d derivations, and \mathfrak{d} is the set of dimension = d derivations. To efficiently determine the values for all = d+1 derivations, newton solves a *linearized* program.

We discuss solving the linear subproblems in $\S 5.3$. Unfortunately, not all subproblems \mathfrak{d}' can be solved exactly *and* efficiently for all semirings; important

⁹⁹Our pseudocode generalizes Schlund (2016, Definition 3.9), which is for the special case of context-free grammars, to programs. The generalization is relatively seamless thanks to the building blocks from §3.1.

Algorithm 11 newton: Solve program p by Newton's method

```
1. def newton(p):
          \langle \mathfrak{m}, \mathfrak{d} \rangle \leftarrow \text{newton\_init}(\mathfrak{p})
          while true:
               \langle \mathfrak{m}', \mathfrak{d}' \rangle \leftarrow \text{newton step}(\mathfrak{p}, \mathfrak{m}, \mathfrak{d})
              if \mathfrak{m}' = \mathfrak{m}: break
               \langle \mathfrak{m}, \mathfrak{d} \rangle \leftarrow \langle \mathfrak{m}', \mathfrak{d}' \rangle
          return m
 8. def newton step(\mathfrak{p}, \mathfrak{m}, \mathfrak{d}):
         \mathfrak{m}' \leftarrow \mathfrak{m} \oplus \mathfrak{d}
                                          Should omit constants and builtins; these are added by line 11.2.
         ð' ← ? \ <
         for (x \oplus \exists y_1 \circ \cdots \circ y_K) \in \mathfrak{p}:
              \triangleright Base case: combine at least two derivations with dimension =d
              for 1 \le i < j \le K:
                   \mathfrak{d}'[x] \oplus \mathfrak{m}'[y_{[:i)}] \circ \mathfrak{d}[y_i] \circ \mathfrak{m}[y_{(i:j)}] \circ \mathfrak{d}[y_j] \circ \mathfrak{m}[y_{(i:l)}]
              \triangleright Recursive case: derive = d+1 derivations by combining one = d+1 derivation
                  with \leq d \ derivations
              for k \in [1:K] where y_k \notin \mathbb{V} \cup \mathfrak{B}:
                    \mathfrak{d}'[x] \oplus = \mathfrak{m}'[y_{[:k)}] \circ y_k \circ \mathfrak{m}'[y_{(k:]}]
          return \langle \mathfrak{m}', \text{ solve linear}(\mathfrak{d}') \rangle
                                                                                          > Solve using specialized linear solver (§5.3)
19. def newton init(\mathfrak{p}):
          ▷ Initial set of new derivations
         \mathfrak{m} \leftarrow \langle \langle \rangle; \mathfrak{d}' \leftarrow \langle \rangle \rangle
         for (x \oplus \exists y_1 \circ \cdots \circ y_K) \in \mathfrak{p} where at most 1 of y_1, \ldots, y_K \notin \mathbb{V} \cup \mathfrak{B}:
              \mathfrak{d}'[x] \oplus = y_1 \circ \cdots \circ y_K
         return \langle \mathfrak{m}, \text{solve linear}(\mathfrak{d}') \rangle
                                                                                             ▷ Solve using specialized linear solver (§5.3)
```

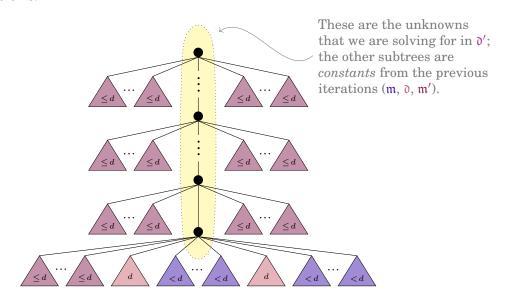
tractable special cases include commutative semirings and cases where the subproblems are either left linear or right linear. Other cases are discussed in the literature (e.g., Esparza et al. (2007), Schlund (2016), and Reps et al. (2017)). We note that fixpoint iteration methods are always available for (approximately) solving \mathfrak{d}' ; however, fixpoint iteration may not converge quickly (or at all).

As noted in other work on Newton's method (e.g., Schlund (2016)), decomposing the program into strongly connected components—as we did for prioritization (§5.1.4)—can speed the algorithm up considerably. Our hybrid algorithm (§5.4) speeds up Newton's method in this way.

We argue for the correctness of newton in the following proposition.

Proposition 4. On iteration d of Algorithm 11 (line 11.6), the value of \mathfrak{m} is the valuation function restricted to dimensions < d, and the value of \mathfrak{d} is the valuation function restricted to dimension = d.

Proof (Sketch). We first argue that the value \mathfrak{d}' returned by the call to newton_step on line 11.4 is precisely the sum of dimension =d+1 derivations. Below, we illustrate how to build dimension =d+1 derivations inductively from the set of dimension =d derivations, the set of dimension $\leq d$ derivations, and the set of $\leq d$ derivations.



We argue that our scheme for creating new derivations generates derivations of dimension =d+1. According to Definition 29, to create a dimension $\geq d+1$ derivation, we need to combine at least two =d derivations. That is why the base case of newton_step (lines 11.13–11.14) combines at least two =d derivations;

this is depicted as the rule connecting at least two adjacent =d derivations; this step creates a subderivation with dimension equal to =d+1. Now, because all subderivations that can combine along the spine (the path up to the root of the derivations) of the new derivation have dimension $\leq d$, the new derivation has dimension =d+1; we see this in the recursive case of newton_step (lines 11.15–11.17). Thus, the dimension of the new derivation is precisely dimension =d+1. To ensure that the enumeration of =d+1 is exhaustive, the program \mathfrak{d}' includes every rule that can combine the $\leq d$ subderivations and the base case of at least two =d subderivations. The program \mathfrak{d}' is a compact encoding of the value of this potentially infinite set of =d+1 derivations. Its solution is found by calling a specialized solver (described in §5.3) for programs with linear rules.

On line 11.9, we see that \mathfrak{m}' is the sum of \mathfrak{d} and \mathfrak{d} . On line 11.6, \mathfrak{m}' overwrites \mathfrak{m} . Thus, at iteration d, \mathfrak{m} the total value of all derivations with < d.

5.3 Specialized Right-Linear Solver

This section gives an algorithm that solves the special family of right-linear programs, such as those generated by Newton's method (§5.2).¹⁰⁰ Unlike fixpoint iteration and Newton's algorithm, this algorithm is exact: it can sum over the infinite set of right-linear derivations in closed form. Until now, our only exact algorithms were for acyclic programs (i.e., those defining only finitely many derivations). We can perform inference efficiently because the solution can be

¹⁰⁰We note that all algorithms we derive in this section for the right-linear case work equally well for the left-linear case with straightforward modifications.

expressed in terms of finitely many regular operations. A simple example of a right-linear program is

```
127 x += 1.
128 x += 0.5 * x.
129 a += 0.5.
```

Although this program has infinitely many derivations of x, we can solve for the total value of x using grade school algebra (see earlier discussion Example 10). Right-linearity can be checked syntactically: each rule has one of two possible forms (1) a completely constant body side like #127, or (2) the body contains a single subgoal as its rightmost subgoal like #128 (that subgoal cannot be a builtin). A program can often be put into a right-linear form using simple program transformations. The example below is not (immediately) right-linear.

However, if we unfold the a subgoal in #131 and exploit the commutativity of multiplication (if permitted by the program's semiring) to reorder the right-hand side, then we arrive at the right-linear program from before. Unfortunately, not all semirings are commutative; thus, permuting the right-hand side is not allowed. The following example, using the semiring of formal languages (§3.3), is linear but not right-linear.¹⁰¹

Unlike the previous example, we cannot transform this program into a rightlinear program because its semiring is not commutative. Unfortunately, the

¹⁰¹The solution to this program is $\mathbf{x} = \{a^n b^n \mid n \in \mathbb{N}\}$, which is not a regular language. Thus, it cannot be expressed as a (finite) expression involving the operations \circ , \oplus , and $(\cdot)^*$.

right-linear algorithm cannot solve this program efficiently; we would need a general linear algorithm.

Algorithm 12 leverages a general-purpose algorithm for solving right-linear recurrences, i.e., equations of the form $\mathbf{x} = \mathbf{b} \oplus \mathbf{A} \mathbf{x}$ where $\mathbf{A} \in \mathbb{V}^{N \times N}, \mathbf{b} \in \mathbb{V}^{N}$. The solution is found by the matrix star operation $\mathbf{x} = \mathbf{A}^* \mathbf{b}$, which we defined in Definition 21. Special cases of this algorithm for specific applications are well-known, e.g.,

- matrix inversion (Gauss–Jordan Elimination)¹⁰²
- all-pairs shortest path (Floyd, 1962; Warshall, 1962; Roy, 1959)
- finite-state automaton to regular expressions conversion (Kleene, 1956)
- transitive closure (Warshall, 1962)

These problems were unified into what is often called **algebraic path problems** (Lehmann, 1977; Tarjan, 1981). Unifying these problems was just the (closed) semiring abstraction (§3.3) at work again! A choice of semiring—with the same underlying algorithm that sums over all paths in a graph—solves all the above problems and many others.

A useful perspective on the matrix star operation is that it computes the total weight of paths in a weighted directed graph. Paths may be regarded as a restricted family of right-recursive derivations (Definition 23). Suppose the

¹⁰²Althoen and McLaughlin (1987) provide an interesting historical overview of the Gauss–Jordan Elimination algorithm.

¹⁰³That is the interpretation used in Algorithm 12. Similarly, Algorithm 11 uses this interpretation to sum over linearly recursive derivations. Notice, however, that Algorithm 11 requires solutions to be linearly recursive. To use the efficient right-recursive solver, we need the subgoals to be permutable (e.g., a commutative semiring).

Algorithm 12 solve_linear: Solve a right-linear program p by reduction to the matrix star operation (Definition 21).

```
1. def solve linear(\mathfrak{p}):
        \overline{\mathfrak{m}} \leftarrow \overline{\operatorname{support}}(\mathfrak{p})
                                                                                         ▷ Estimate p's support (Algorithm 9)
        N \leftarrow \{ \mathbf{x} \mid \overline{\mathfrak{m}}[\mathbf{x}] = \mathsf{true} \}
        \mathbf{A} \leftarrow \mathbf{0}^{N \times N}; \mathbf{b} \leftarrow \mathbf{0}^{N}
        ▷ The following two loops instantiate p's rules against the estimated support. In
           the loop's body, we assemble the matrix A and vector b that will be passed to the
           right-linear solver.
        for (x \oplus = y_1 \circ \cdots \circ y_K) \in \mathfrak{p}:
           for \underline{\phantom{a}} \in \overline{\mathfrak{m}}[y_1 \circ \cdots \circ y_K] :
                if x \oplus = y_1 \circ \cdots \circ y_K is not ground:
                    raise Error("nonground reasoning not supported")
                else if y_1 \circ \cdots \circ y_K = v where v \in \mathbb{V}:
                    \mathbf{b}_{r} \oplus = \mathbf{v}
                else if y_1 \circ \cdots \circ y_K = v \circ y where v \in \mathbb{V}, \ y \notin \mathfrak{B}^{104} :
                else
                    raise Error("program is not right linear")
15
        > Solve the right-linear system with the matrix star operation (Definition 21).
        \mathbf{s} \leftarrow \mathbf{A}^* \mathbf{b}
        ▷ Convert solution vector s into a program in the obvious way
        \mathfrak{m} \leftarrow 75
        for x \in N : \mathfrak{m}[x] \oplus = \mathfrak{s}_x
20
        return m
```

matrix A is the weighted adjacency matrix of a graph. Then A_{ik} is the total weight of all paths of length =1. Similarly, the total weight of all paths between i and k of length =2 is $(\mathbf{A} \mathbf{A})_{ik}$. Thus, the star operation on matrices $\mathbf{A}^* \stackrel{\text{def}}{=} \sum_{i=0}^{\infty} \mathbf{A}^i$ computes the total weight of paths of length $0, 1, \ldots$ and adds them up.

Definition 21 provides a simple recursive method for computing A^* . It runs in $\mathcal{O}(|N|^3)$ time and can be made to run in $\mathcal{O}(|N|^2)$ space.¹⁰⁵ The runtime can be improved by decomposition into strongly connected components (Definition 12).¹⁰⁶ Our hybrid algorithm (§5.4) performs this optimization, which is

 $^{^{104}}$ We may take v=1 to match the pattern.

¹⁰⁵Sub-cubic algorithms exist when V is a *ring*, e.g., Strassen (1969).

¹⁰⁶The SCC-based algorithms for transitive closure date back to Purdom (1970).

not in Algorithm 12. Note that there is a speed–accuracy tradeoff when using the right-linear solver. Often fixpoint iteration can reach a tolerable level of accuracy faster than the exact algorithm. Lastly, we mention that solve_linear requires the program's rules to come to ground (see line 12.8). 108

5.4 Hybrid Algorithm

This chapter has provided several algorithms that each work best under different conditions. This section provides a hybrid algorithm that tries to use the best algorithms available to solve a given program \mathfrak{p} . We saw that if \mathfrak{p} is acyclic, it can be solved efficiently by agenda with toposorted prioritization (§5.1.4). However, if \mathfrak{p} is cyclical, it is best to solve \mathfrak{p} by breaking it up into cyclical subprograms $\mathfrak{p}^{(1)}, \ldots, \mathfrak{p}^{(M)}$ by strongly connected component (SCC) decomposition. ¹⁰⁹ This decomposition allows each subprogram to be solved separately, provided that the solutions to the subprograms it depends on have been computed beforehand; this can be achieved by solving the subprograms in a topological order. By analyzing the structure of $\mathfrak{p}^{(m)}$, we can determine what the best method for

¹⁰⁷The convergence rate of fixpoint iteration depends on the properties of the A matrix. In the case of the real semiring, we can bound the number of iterations to reach a given error tolerance in terms of the spectral radius of A.

 $^{^{108}}$ It is often possible to work around this restriction. For example, the abbreviation transformation (§6.2) can remove certain nonground rules that cause the assertion to fail. For example, the program $\{x(I) += 1. x(I) += 0.5 * x(I).\}$ is linear, but unfortunately not ground because of the free variable I. Abbreviation transforms it to $\{x' += 1. x' += 0.5 * x'.\}$, which can be solved by solve_linear. It also provides a recovery rule x(I) += x' to rebuild the original nonground items.

¹⁰⁹SCC decomposition is a common optimization for solving cyclical equations (e.g., Goodman (1998), Nederhof and Satta (2008), Esparza et al. (2007), Schlund et al. (2013), Mohri (2002), Eisner (2023), Baharev et al. (2016), Chiang, McDonald, et al. (2023), and Purdom (1970)).

solving it is likely to be. Below is our ranking:

- 1. solve_linear: Because it is exact. However, it requires a closed semiring (§5.3), at least one of {right-linearity, left-linearity¹⁰⁰, linearity+commutativity}, and that the program can be grounded (see line 12.8).
- newton: Because it converges more quickly than fixpoint iteration methods.
 However, it requires that the induced linear subproblems can be solved efficiently.
- 3. agenda: When all else fails, the subprogram can be solved with agenda, as it runs under the most general conditions.

Algorithm 13 provides pseudocode for our hybrid algorithm, which solves each subprogram in a topological order where the solutions to earlier subprograms are taken as inputs to later subprograms. It makes use of a subroutine best_solve that uses the ranking of solvers above to select the best strategy for solving the given subprogram.¹¹⁰

Discussion.

- The discussion related to the termination of support applies here (§5.1.4).
- When solving an SCC with either newton or a fixpoint algorithm, it may be necessary to use a numerical convergence threshold (or iteration limit), as

¹¹⁰ A simple implementation of best_solve would use exception handling to move down the list of preferred solvers. This can be sped up when the necessary properties can be checked ahead of any expensive computations. For example, linearity can be checked without instantiating the program's rules in Algorithm 12.

Algorithm 13 solve: Solve program p using the hybrid algorithm.

```
1. def solve(p):
        \overline{\mathfrak{m}} \leftarrow \overline{\operatorname{support}}(\mathfrak{p})
                                                                                           ▷ Estimate p's support (Algorithm 9)
        \mathfrak{m} \leftarrow 75
        for \sigma \in \operatorname{sccs}(\operatorname{incoming}, \overline{\mathfrak{m}}):
                                                                     ▷ Topologically ordered list of SCCs (Algorithm 5)
            \triangleright Restrict \mathfrak{p} to compute the items in the current component, \sigma.
            \mathfrak{p}' \leftarrow ? 
            for x \in \sigma:
                for (y_1 \circ \cdots \circ y_K) \in \mathfrak{p}[x]:
                                                                ⊳ Add the specialized rules to the subprogram;
                                                                   Note: p's rules may participate in multiple SCCs.
                    \mathfrak{p}'[x] \oplus = y_1 \circ \cdots \circ y_K
            \triangleright Solve the subprogram \mathfrak{p}' with the best applicable method, taking the solutions to
               earlier subprograms as input through m. 111
            \mathfrak{m} \leftarrow \mathsf{best} \; \mathsf{solve}(\mathfrak{p}' \oplus \mathfrak{m})
11.
       return m
```

discussed earlier in this chapter. Otherwise, the Algorithm 13 could run forever on a single subprogram with making progress on the others.

• Looking forward to other chapters, we note that it is often best to *optimize* the program before executing. We note, in particular, the program specialization techniques of (§6.2) and the program optimizer (§8). These methods produce a semantically equivalent program that is more efficient to execute.

Summary

In this chapter, we described many algorithms for solving Dyna programs. We gave three implementations of fixpoint iteration. We extended Eisner, Goldlust, et al. (2005)'s forward chaining algorithm to allow nonground updates and delayed constraints. We adapted Esparza et al. (2007)'s semiring-weighted Newton algorithm to Dyna programs, which provides a powerful new solver to 111 The ⊕-addition below should be implemented carefully to avoid the cost of copying.

the Dyna toolbox. We also gave a specialized solver for programs with a certain linear structure. Lastly, we gave a unified, hybrid algorithm that breaks up a large program into smaller subprograms that can be solved separately with the best applicable method.

Chapter 6

Program Analysis

This chapter is about static program analysis: proving that the program satisfies certain properties without actually running it on *all possible* input data sets. Analyzing program behavior without actually running the program on concrete data is known in the programming languages community as *abstract interpretation* (Cousot and Cousot, 1977); it can be viewed as the partial evaluation of a program to acquire constraints on its *possible* semantics.

• Abstract analyses are useful for code generation because code generation can take advantage of any guarantee that an item will have a particular type or property. In §6.2, we leverage the output of abstract analysis to generate specializations of the program (i.e., Dyna-to-Dyna code generation) that can compute a subset of the items more efficiently by removing unnecessary intermediate items.

¹¹²In the context to code generation into an object-oriented host language, these guarantees can be used to represent the item efficiently as an instance of a class specialized to objects of that type and to speed up tests (e.g., pattern-matching).

• Abstract analyses are useful for complexity analysis: If the program stores and iterates over a collection of distinct objects of a particular type, then knowing the type's cardinality on the concrete data can yield bounds on the program's space and runtime requirements (§6.3); our program optimizer (§8) uses such bounds as cost function.

Below is a listing of this chapter's contents.

6.1	Type A	Analysis
	6.1.1	Formulation
	6.1.2	Abstract Forward Chaining
6.2	Type-	based Program Specialization
	6.2.1	Inferring a Type for Useful Items
	6.2.2	Specialization Transformation
	6.2.3	Abbreviation Transformation
6.3	Auton	natic Space and Time Analysis
	6.3.1	Space Analysis
	6.3.2	Time Analysis
6.4	Applio	cations

6.1 Type Analysis

6.1.1 Formulation

This section describes a system that infers the set of items that may have a non- $\underline{0}$ value in a program $\mathfrak p$ for all input data $\mathfrak D$ that satisfy user-provided input constraints $\overline{\mathfrak D}$. We call the set of non- $\underline{0}$ -valued items the **support** of the program. Let $\mathsf{nz}(\mathfrak D)$ and $\mathsf{nz}(\mathfrak p(\mathfrak D))$ denote the support of their respective arguments. Since we do not know the *concrete* data $\mathfrak D$ to execute $\mathfrak p(\mathfrak D)$, we compute an upper bound $\overline{\mathfrak m}$ on its support through *abstract* interpretation such that 114

$$\underbrace{nz(\mathfrak{D})\subseteq\overline{\mathfrak{D}}}_{\mathit{input satisfies constraints}} \Longrightarrow \underbrace{nz(\mathfrak{p}(\mathfrak{D}))\subseteq\overline{\mathfrak{m}}}_{\mathit{valid upper bound on support}}$$

In other words, $\overline{\mathfrak{D}}$ characterizes the subset of \mathbb{H} for which the input items \mathfrak{D} is non- $\underline{0}$, and $\overline{\mathbb{m}}$ characterizes the subset of \mathbb{H} for which the derived items $\mathfrak{p}(\mathfrak{D})$ is non- $\underline{0}$. Therefore, we call $\overline{\mathfrak{D}}$ the **input type** and $\overline{\mathbb{m}}$ the **derived type**.

We now describe our type specification language and how it is used to reason about possible items. Continuing with the CKY example (Example 6), we provide the specification of the input data. However, as we are interested in the *possibility* that an item may have a non-0 value, rather than finding its specific value, we first transform the program as illustrated below.

Example 18. Recall, the program to analyze is

 $[\]overline{}^{113}$ Recall that $\mathfrak{p}(\mathfrak{D})$ is shorthand for $\mathfrak{p}'[\cdot]$ and $\mathfrak{p}'=\mathfrak{p}\oplus\mathfrak{D}$.

¹¹⁴For naturally ordered semirings (Definition 26), it should be possible to refine our technique to produce an upper bound on the *value* of each derived item as a function of upper bounds on the values of its input items.

```
params: word; len; \gamma.

\beta(X,I,K) += \gamma(X,Y,Z) * \beta(Y,I,J) * \beta(Z,J,K).

\beta(X,I,K) += \gamma(X,Y) * \beta(Y,I,K).

\beta(X,I,K) += \gamma(X,W) * word(W,I,K).

\beta(X,I,K) += \beta(X,W) * len(N).
```

Because no rules have been specified to define the word, len, and γ items, one might think that these items have value 0.115 However, the params declaration says that the program's input will define these items.

The CKY program $\mathfrak p$ above can be automatically "booleanized" to yield the **type program** $\overline{\mathfrak p}$ below. 116 The :- and , symbols traditionally denote the \oplus = and \circ operations in the Boolean semiring.

```
140 params: \overline{\text{word}}; \overline{\text{len}}; \overline{\gamma}.

141 \overline{\beta}(X,I,K) := \overline{\gamma}(X,Y,Z), \overline{\beta}(Y,I,J), \overline{\beta}(Z,J,K).

142 \overline{\beta}(X,I,K) := \overline{\gamma}(X,Y), \overline{\beta}(Y,I,K).

143 \overline{\beta}(X,I,K) := \overline{\gamma}(X,W), \overline{\text{word}}(W,I,K).

144 \overline{\text{goal}} := \overline{\beta}(S,0,N), \overline{\text{len}}(N).
```

The intention is that $\overline{\beta}(X,I,K)$ will have value true in the type program for all X,I,K such that $\beta(X,I,K)$ has a non- $\underline{0}$ value in the CKY program. Of course, each program needs input: if \mathfrak{p} is run on data \mathfrak{D} , then $\overline{\mathfrak{p}}$ must be run on a booleanized version $\overline{\mathfrak{D}}$. The valuation $\overline{\mathfrak{p}}(\overline{\mathfrak{D}})$ now effectively specifies an inferred type of items that *might* be derived by the program.¹¹⁷

To abstract away from the specific input data \mathfrak{D} , users may *directly* specify the booleanized input data $\overline{\mathfrak{D}}$ (mentioned above). This $\overline{\mathfrak{D}}$ may be any upper bound on its support. $\overline{\mathfrak{D}}$ serves as a type declaration for the unknown input \mathfrak{D} .

¹¹⁵In other words, a program that is expecting inputs is only *partially specified*. When input items are queried, they are returned as delayed constraints (rather than incorrectly assigning them a value of 0).

¹¹⁶The translation additionally maps any constant values to false or true.

¹¹⁷For example, $\overline{\mathfrak{p}}(\overline{\mathfrak{D}})[\![\overline{\beta}(\mathsf{np},3,5)]\!] = \mathsf{true} \text{ if } \mathfrak{p}(\mathfrak{D})[\![\beta(\mathsf{np},3,5)]\!] \neq \underline{0}$ —and possibly even if $\mathfrak{p}(\mathfrak{D})[\![\beta(\mathsf{np},3,5)]\!] = \underline{0}$, for example, if $\underline{0}$ can result from adding non- $\underline{0}$ values under \oplus . Such overestimates can arise because the boolean program does not track actual values.

The type defined by $\overline{\mathfrak{p}}(\overline{\mathfrak{D}})$ then continues to include all items in $nz(\mathfrak{p}(\mathfrak{D}))$, for any and every dataset \mathfrak{D} for which this upper bound holds. For example, $\overline{\mathfrak{p}}(\overline{\mathfrak{D}})$ still gives value true to $\overline{\beta}(np,3,5)$ if $\beta(np,3,5)$ could possibly have a non- $\underline{0}$ value.

An example of an input type declaration for CKY is the following boolean program $\overline{\mathfrak{D}}$:

```
145  params: k; w; n.

146  \overline{\text{word}}(\text{W:w,I:n,K:n}) := \text{I} < \text{K}.

147  \overline{\text{len}}(\text{N:n}).

148  \overline{\gamma}(\text{X:k,Y:k,Z:k}).

149  \overline{\gamma}(\text{X:k,Y:k}).

150  \overline{\gamma}(\text{X:k,W:w}).
```

The lines #146–150 specify 5 sets whose union is a type that contains all input items. We use the notation $\overline{\text{word}}(W:w,I:n,K:n) := I < K$ in #146 as a shorthand for $\overline{\text{word}}(W,I,K) := w(W)$, n(I), n(K), I < K. This type indicates that items in the set $\{\text{word}(W,I,K) \mid w(W),n(I),n(K),I < K\}$ may have non-0 values in the input data $\mathfrak{D}.^{118,119}$ The predicates such as n could be defined by additional lines like $n(I) := 0 \le I \le 20$, which defines the possible sentence positions for sentences of length 20. However, to abstract away from any particular sentence length, we omit such a definition, leaving n as a **parameter** (#145). Thus, the above program $\overline{\mathfrak{D}}$ actually defines a *parametric* type, just as List[α] in many programming languages is a parametric type that depends on a concrete choice such as

¹¹⁸Note that this input type is compatible with *lattice parsing* (e.g., Tomita (1986) and Hall (2005)); in that setting, the nodes of the lattice (i.e., weighted acyclic finite-state automaton) would be labeled by their position in a topological ordering of the input, and word(W, I, K) encodes the weight of the arc $I \xrightarrow{W} K$.

¹¹⁹A natural alternative to #146 is word(W:w,I:n,K:n) :- (K is I + 1), which uses the built-in addition constraint. This type denotes the set {word(W,I,I+1) | w(W),n(I),n(I+1)} and it directly encodes the description of the CKY input data that we gave in Example 6. Although this input characterization may be more precise for some applications, the conclusions that we can derive from it appear to be the same (i.e., we are unable to infer a tighter derived type).

 $\alpha = \text{string.}$ It depends on three other relations, which specify the allowed terminals, nonterminals, and sentence positions. The concrete input type would have the form $\overline{\mathfrak{D}}(k,w,n)$. Our program analyzer reasons about these user parameters symbolically; thus, it performs *symbolic evaluation* (e.g., Baldoni et al. (2018)).

Our abstract forward chaining algorithm (§6.1.2) will derive the following rules that specify a valid type for the derived items of $\overline{\mathfrak{p}}(\overline{\mathfrak{D}}(k,w,n))$:

```
\overline{\beta}(X:k,I:n,K:n) := I < K.
\overline{goal}.
```

Adding these rules and $\overline{\mathfrak{D}}$ yields the desired description of a type that covers all non- $\underline{0}$ items. Call the resulting program $\overline{\mathfrak{m}}$. The program $\overline{\mathfrak{m}}$ uses the same parameters as $\overline{\mathfrak{D}}$ and provides an upper bound on $\overline{\mathfrak{p}}(\overline{\mathfrak{D}}(k,w,n))$. The bound implies that we will only build $\overline{\mathfrak{p}}$ items where X is a nonterminal (has type k), and the I and K variables are sentence positions (have type n) with I < K. Additionally, the bound says (reassuringly) that it may be possible to build a goal item. If goal were impossible, we would have detected a mistake in the code. Similarly, if the bound had included $\overline{\mathfrak{p}}(X:k,I:k,K:n)$, it would be a suspiciously loose upper bound, since I should never be a nonterminal.

Notice that $\overline{\mathfrak{p}}' = \overline{\mathfrak{p}} \oplus \overline{\mathfrak{D}}$ would *also* be a valid type for $\overline{\mathfrak{p}}(\overline{\mathfrak{D}}(k,w,n))$. However, $\overline{\mathfrak{p}}'$ is not easily consumed by downstream applications. Our analyzer (§6.1.2) is guaranteed to output a derived type, which is a program (like $\overline{\mathfrak{m}}$ above) where each rule is *simple*. We define a rule to be **simple** if its body contains only constants, user parameters, and built-in constraints. Simple rules may be regarded as the *symbolic* generalization of *axioms* (§2).

¹²⁰ E.g., space and time complexity analysis (§6.3) and program specialization (§6.2).

In summary, our type analyzer is given the program $\mathfrak p$ together with a parametric input type $\overline{\mathfrak D}$, and computes a parametric derived type such as #151-152 for the items derived by $\mathfrak p$. In effect, it runs the booleanized program $\overline{\mathfrak p}$ on the input upper bound $\overline{\mathfrak D}$. However, as the derived type is only required to be an upper bound, we will permit the analyzer to "round up" as it goes, as needed to keep the type simple. 121

6.1.2 Abstract Forward Chaining

We now develop our type inference algorithm, which is based on fixpoint iteration (§4.3) in the boolean semiring. Here, we generalize fixpoint iteration's valuation estimate ν to an abstract valuation where the user's type parameters (such as n) are used to define ν symbolically, as we saw in the derived type inferred for the CKY example. This generalization is necessary because we no longer have a concrete upper bound on the input data with which to initialize forward chaining: the n, k, and w parameters are unknown relations. Instead, we specify parametric rules that upper-bound the unknown input data. Forward chaining derives new simple rules, which may also refer to the relations like n and <. Indeed, these relations may emerge in the derived type \overline{m} , as we saw in #151.

We provide the abstract forward chaining algorithm in §6.1.2.1. We address the challenges arising from constraint accumulation in §§6.1.2.2–6.1.2.5.

¹²¹The abstract interpretation literature often calls rounding up *widening*.

¹²²Recall that fixpoint iteration is also the basis for ordinary forward chaining (Ullman, 1988; Eisner, Goldlust, et al., 2005) and many of the algorithms that we developed in §5.

6.1.2.1 Type Inference Algorithm

To generalize forward chaining in the boolean semiring, we will represent the valuation ν from §4 as a program \overline{m} —just like we did in §5. We will reason in the boolean semiring. This representation is helpful because we can use programs to represent unknown and infinite sets. We require each rule in \overline{m} to have a specific syntactic form: each must be a $simple\ type$. A $simple\ type$ is a rule of the form $\overline{x}:-c_1,\ldots,c_M$ where each c_m is externally defined as either a built-in constraint or a parametric constraint and may contain variables. Note that \overline{x} cannot appear in the body of any other simple type. We saw the derived type for the CKY example #146-152 is the collection of simple types.

In some sense, program rules #141–150 already describe the derived type, but they do so indirectly and not as a collection of simple types. This form ensures that each simple type is self-contained and can be implemented as a (nonrecursive) object-oriented class.

Type inference is performed by fixpoint iteration on a symbolic step operation. This should be viewed as a "lifted" version of the forward chaining algorithm from §5: it derives simple types from other simple types rather than deriving ground items from other ground items. The constraints that appear in these simple types are passed around during execution. Such constraints are known in the constraint logic programming literature as "delayed constraints," in contrast to constraints like 2 < 3 or 3 < 2 that can be immediately evaluated and replaced with true or false (respectively).

Algorithm 14 abstract_fc: Abstract forward chaining algorithm for type program \overline{p} . This procedure assumes that the input type has been included in \overline{p} .

However, the challenge that arises from the repeated application of step is that constraints may accumulate. They can accumulate indefinitely in the case of recursive rules since they effectively unroll an unbounded chain of recursive calls, as will be illustrated in §6.1.2.2. To prevent divergence, we propose the following strategies: constraint relaxation (relax; §6.1.2.2), and constraint propagation (propagate; §6.1.2.3). We discuss termination in §6.1.2.5. We give pseudocode for the overall approach in Algorithm 14.

6.1.2.2 Relaxation

Consider the following input type $\overline{\mathfrak{D}}$ for Example 4.

```
params: n.
stop(S) :- n(S).
edge(S,S') :- n(S), n(S').
```

Applying the step operation once, which renames variables to avoid conflicts, we

¹²³Note that when user parameters are queried in $\overline{\mathfrak{p}}$, they are returned as delayed constraints by Algorithm 2. For example, let $\overline{\mathfrak{p}}$ be the booleanized CKY program together with its input type. Then, if we perform the lookup $\overline{\mathfrak{p}}[n(K)]$, we would get the result stream containing n(K) as its only element.

 $[\]overline{m}$ as a set of rules. In our implementation, this set replaces rules with their most general version, i.e., if there exists a pair of simple types r, r' such that $r' \subseteq r$, we remove r' and keep r. This ensures, among other things, that rules that are equivalent under variable renaming do not both appear in the set. It can also help to achieve convergence. See §6.1.2.4 for further discussion.

derive the following simple type for β items (other types not shown):

```
\overline{\beta}(S_1) := n(S_1).
```

Applying step a second time, we derive

```
157 \overline{\beta}(S_1) :- n(S_1).
158 \overline{\beta}(S_1) :- n(S_1), n(S_2).
```

Notice in the second rule that $n(S_2)$ is lingering as a constraint from the recursive call. In keeping with the (Prolog-like) semantics of boolean programs, the variable S_2 is existentially quantified since it is a "local" variable that appears only in the body of the (:-) rule. After applying step k times, we derive

```
159 \overline{\beta}(S_1) := n(S_1).

160 \overline{\beta}(S_1) := n(S_1), n(S_2).

161 \vdots

162 \overline{\beta}(S_1) := n(S_1), n(S_2), n(S_3), ..., n(S_k).
```

This procedure is diverging by elaborating the type to describe paths of every length, rather than converging to a fixpoint (i.e., a type that describes all nodes that can reach a terminal state on paths of *any* length). Our strategy for eliminating constraint accumulation from previous levels of recursion is to drop (i.e., relax) any constraint that refers to local variables:

Definition 30 (relax). Let \bar{x} :- c_1 , ..., c_M be a simple type. The operation relax drops each delayed constraint c_m that refers to a variable that does not appear in \bar{x} . In other words, we obtain a new simple type with the same head, but with the following subset of constraints $\{c_m \mid \text{vars}(c_m) \subseteq \text{vars}(\bar{x}), m \in [1:M]\}$. To relax a set of simple types, we simply map relax over the set.

Applying relax to the second application of step to #157–158, we see that it would drop the second constraint of the second simple type because it depends

on the non-head variable S2, yielding:

```
\frac{\overline{\beta}}{\beta}(S_1) := n(S_1).

\frac{\overline{\beta}}{\beta}(S_1) := n(S_1).
```

We may merge these rules because they are duplicates. We now have the following type

```
\overline{\beta}(S_1) := n(S_1).
```

We now have a fixpoint: this iteration's type (after relaxation) is equal to the previous iteration's type (#156). Notice that some precision is lost because $\overline{\beta}(S)$ is true even for nodes S for which there does not exist a path ending in \overline{Stop} .

Why is relax valid? In general, dropping constraints from a simple type simply makes the type larger, raising the upper bound. Dropping constraints that depend on non-head variables is simply a convenient choice.

We now consider an example where relax leads to an overly loose type. Recall CKY (Example 18). The simple types from $\overline{\mathfrak{D}}$ (#145–150) are added to $\overline{\mathfrak{m}}$ on the first step, such as

```
166 \overline{\gamma}(X,W) := k(X), w(W).
167 \overline{w} word (W,I,K) := w(W), n(I), n(K), I < K.

On the next step, \#143 from \overline{p}
168 \overline{\beta}(X,I,K) := \overline{\gamma}(X,W), \overline{w} word (W,I,K).

gives us an initial type for \beta:
169 \overline{\beta}(X,I,K) := k(X), w(W), n(I), n(K), I < K.

After relax, we have
170 \overline{\beta}(X,I,K) := k(X), n(I), n(K), I < K.

However, on the next step, \#141 from \overline{p}
171 \overline{\beta}(X,I,K) := \overline{\gamma}(X,Y,Z), \overline{\beta}(Y,I,J), \overline{\beta}(Z,J,K).
```

will expand $\overline{\mathfrak{m}}$ to the following type:

```
\overline{\beta}(X,I,K) := k(X), k(Y), k(Z), n(I), n(J), n(K), I < J, J < K.
```

relax will drop all the constraints on local variables:

```
\overline{\beta}(X,I,K) := k(X), n(I), n(K).
```

Unfortunately, relaxation resulted in the ordering information (<) being lost. In the next section, we will see how to combine I < J and J < K (before relaxing them) to deduce that I < K. Since I < K only constrains head variables, it will survive the relaxation step—allowing the analyzer to deduce that $\overline{\beta}(X,I,K)$ items all have property I < K (rather than the first generation having I < K, the second having ($\exists J$) I < J, J < K, the third having ($\exists J_1, J_2, J_3$) I < $J_1, J_1 < J_2, J_2 < J_3, J_3 < K$, and so on.).

6.1.2.3 Constraint Propagation

Constraint propagation rules (Frühwirth, 1998), are commonly used in constraint logic programming, such as ECLiPSe (Apt and Wallace, 2007), computer algebra systems, such as Mathematica (Wolfram, 2003), and theorem provers, such as Z3 (Moura and Bjørner, 2008). In our setting, they provide a mechanism for users to supply domain-specific knowledge about type parameters and built-in constraints. They are used to derive new constraints during inference, which may serve to counteract excessive relaxation (§6.1.2.2).

Recall the loose type inferred in the CKY example ($^{\#}172$; repeated below):

```
\overline{\beta}(X,I,K) := k(X), k(Y), k(Z), n(I), n(J), n(K), I < J, J < K.
```

We saw that relax causes the ordering relationship I < κ to be lost in #173. Con-

straint propagation rules combat this loss by inferring the additional constraints, such as I < κ by the transitive property of <: 125

```
\overline{\beta}(X,I,K) := k(X), k(Y), k(Z), n(I), n(J), n(K), I < J, J < K, I < K.
```

After applying relax, we arrive at our desired type

```
\overline{\beta}(X, I, K) := k(X), n(I), n(K), I < K.
```

Below, we give motivating examples of constraint propagation rules and how they are useful for our running examples. Constraint propagation rules are specified using the <== syntax seen in the examples below. They may be specified along with the input type.

• In CKY (Example 6), the (constant) grammar start symbol s is an element of the parametric nonterminal type k, and 0 is an element of the parametric sentence position type n:

```
177 	 k(s) \le true.
178 	 n(0) \le true.
```

• To assert that the k and n types are disjoint:

```
179 fail <== k(X), n(X).
```

Here fail is a special constraint that is never satisfied (a contradiction). If k(X) and n(X) ever both appear in the body of a simple type, the simple type is empty and can be deleted. Notice that #177-179 are sufficient to derive a contradiction from either n(s) or k(0) (i.e., both are false).

• To assert that a type d is a subset of a type c:

```
180 \text{ c(X)} \le d(X).
```

¹²⁵The specific syntax for specifying the transitive property is on #181.

¹²⁶We use <== to distinguish propagation rules from program rules and input types.

• To assert that < is a strict partial order:

```
181 (I < K) <== (I < J), (J < K).
182 fail <== (I < I).
```

• To reason about a recursive type, such as a list¹²⁷ of nonterminals (i.e., type k):

```
183 ks([]) <== true.

184 k(X) <== ks([X|Xs]).

185 ks(Xs) <== ks([X|Xs]).
```

Propagation rules for builtins (like those for < above) are provided in a standard library, whereas the user must provide rules for user parameters. If the user does not provide certain rules, our system will still produce a valid upper bound, but it may be looser than necessary. Thus, writing propagation rules may be an interactive experience. However, if the user provides unsound propagation rules, it will likely lead our system to produce an incorrect upper bound.

The propagate method used in abstract forward chaining (abstract_fc) runs a standard constraint propagation algorithm (Frühwirth, 1998)—which is just another forward chaining algorithm—to deduce all transitive consequences of the constraints in the body of a simple type.

Checking constraints. In the propagation step, we check if each of its built-in constraints is satisfied or unsatisfiable. Any satisfied constraint can be safely dropped from its simple type, and it does not loosen the type. If an unsatisfiable constraint exists in a given simple type, then that simple type can be safely dropped from $\overline{\mathbf{m}}$.

¹²⁷The bracket-pipe notation is a shorthand (borrowed from Prolog (Colmerauer and Roussel, 1996)) for working with lists. A list such as [1,2,3] desugars to cons(1,cons(2,cons(3,nil))). The notation [X,Y|Rest] would match as X=1,Y=2,Rest=[3].

6.1.2.4 Removing Redundant Simple Types

In Algorithm 14, we may remove a simple type s from $\overline{\mathfrak{m}}$ if there exists a supertype $t \in \overline{\mathfrak{m}}$ that is $s \subseteq t$ (see also footnote 124). Below is a method for testing such subtype relationships that leverages propagate. We also include a method for merging simple that we use later in the chapter.

Algorithm 15 Subtype, intersection, and merge algorithms for simple types.

```
1. \triangleright Returns true if s \subseteq t; may have false negatives.
2. def subtype(s, t):
       \triangleright (s \subseteq t) \iff (s \cap t) = s
       r \leftarrow \mathsf{intersect}(s,t)
       R \leftarrow \mathsf{propagate}(r)
       S \leftarrow \mathsf{propagate}(s)
       return R = S
                                                                                         ⊳ equal up to variable renaming
8. \triangleright Return a simple type representing the intersection of s and t.
9. def intersect(s, t):
       \theta \leftarrow \text{unify}(\text{head}(s), \text{head}(t))
       s' \leftarrow \boldsymbol{\theta}(s)
       t' \leftarrow \boldsymbol{\theta}(t)
       return (head(s'): - body(s'), body(t'))
14. \triangleright Given simple types r and s, return a simple type t such that r \subseteq t and s \subseteq t. Note that t may
      be a loose bound on r and s.
15. def merge(r,s):
       (r: -\mathcal{R}) \leftarrow r
       (s:-S) \leftarrow s
       X \leftarrow fresh()
                                                    ▶ New variable; used to represent the common head shape
       \mathcal{R}' \leftarrow \text{propagate}(\mathcal{R} \circ (X=r))
       S' \leftarrow \text{propagate}(S \circ (X=s))
       \mathcal{T} \leftarrow \mathcal{R}' \cap \mathcal{S}'
                                                                                         ▷ Only keep common constraints
       return (x : -T)
```

Note that this test may fail to detect the subset relation if the rules for propagate are incomplete. Future systems might use more powerful theorem provers, such as Z3 (Moura and Bjørner, 2008), to test the subset relation.

6.1.2.5 Termination

This section describes a technique for preventing abstract forward chaining from running forever. It ensures finite-time termination with a valid upper bound but may come at the expense of precision. We require a valid upper bound to proceed with downstream analyses, such as the space and time complexity analyzer (§6.3), the program specialization system (§6.2), and the program optimizer (§8).¹²⁸

The challenge in ensuring termination is that abstract forward chaining must provide a valid upper bound, so we cannot simply run it for a finite number of iterations and return the chart because it may underestimate the type. Our approach will be to round up the type (as we did in relax). As an extreme, observe that we can round $\overline{\mathbf{m}}$ all the way up to the unconstrained type X:- true, as it is always a valid upper bound. Unfortunately, this type is completely informative because it is loose for *all* subtypes; having better bounds for at least *some* of the subtypes and looser types for the others would be preferable.

Our termination strategy involves **truncation** according to three usercontrollable parameters $\langle D_{\text{depth}}, D_{\text{builtin}}, D_{\text{propagation}} \rangle$ that work as follows:

terms: We can build infinitely many simple types corresponding to ever-larger terms. To avoid this, we truncate them to an (optional) user-specified depth
 limit D_{depth}. This depth limit will rewrite a simple type in the following way: any subterm that appears in the head at depth D is replaced with a fresh,

They may also be used to ensure that Algorithm 9, which performs support estimation on concrete data, returns a valid upper bound.

unconstrained variable. 129

- **builtins**: The challenge of built-in constraints is that they may introduce an unbounded number of new symbols that can be used in simple types. To prevent this, we keep track of the number of distinct symbols that have been introduced by the builtin relations (across all iterations); once more than D_{builtin} have been introduced, we prevent any of the built-in constraints from evaluating when they attempt to introduce a new symbol.
- propagation: We also must ensure that constraint-propagation rules are not infinitely productive (i.e., they cannot produce an unbounded number of new constraints). To do so, we limit the maximum number of distinct constraints that the constraint propagation rules can produce in a given call to propagate.
 Once more than D_{propagation} constraints have been deduced (in a given simple type), the iterative constraint-propagation procedure stops with whatever it has concluded so for.

These termination heuristics (along with relax) are sufficient to ensure that Algorithm 14 terminates with a valid upper bound. However, when the heuristic is enacted, the resulting analysis will be loose when the heuristic is enacted. Thus, a warning message is provided to the user. If better results are required, the programmer can interact with the system (like an interactive

129 We apply the depth limit modulo structural equality constraints, so that f([1,2,3,4,5]) and f([1,2|Xs]) :- Xs=[3,4,5] are handled the same, as the latter is transformed into the former.

¹³⁰Given finite values for the parameters, only finitely many simple types can be generated. Since the relaxed step operation is monotone nondecreasing, the abstract forward chaining cannot continue to increase $\overline{\mathfrak{m}}$ forever.

theorem prover) to improve the analysis by refining the input type or providing missing propagation rules.

In general, the truncation parameter limits should be set conservatively (to large values); however, there is an efficiency benefit to selecting small values as the type found can be exponentially large in D_{depth} (see Example 19). We cannot simply set the value of D_{depth} to be small in all cases. For example, if $D_{\text{depth}} = 4$, then the simple type $f_1(f_2(f_3(f_4(f_5(X)))) := n(X)$ would be approximated by $f_1(f_2(f_3(f_4(Trunc)))) := true$ which is a much looser type. If the D_{depth} had been ≥ 5 , the analysis would not have lost any precision.

Example 19. Abstract forward chaining produces an exponentially large type (in the D_{depth} parameter) for the program below.

```
bitstring([]).
bitstring([X|Xs]) += bit(X) * bitstring(Xs)
bit(0). bit(1).
```

because it will enumerate exponentially many lists (in D_{depth}) before truncation kicks in. ¹³¹ The example below has $D_{depth} = 5$:

```
 \begin{array}{lll} & \underline{bit(0)}. & \underline{bit(1)}. & \underline{bitstring}([0,0,0]). & \underline{bitstring}([0,0,0,Trunc_1|Trunc_2]). \\ & \underline{bitstring}([]). & \underline{bitstring}([0,0,1]). & \underline{bitstring}([0,0,1,Trunc_1|Trunc_2]). \\ & \underline{bitstring}([1]). & \underline{bitstring}([0,1,0]). & \underline{bitstring}([0,1,0,Trunc_1|Trunc_2]). \\ & \underline{bitstring}([0]). & \underline{bitstring}([0,1,1]). & \underline{bitstring}([0,1,1,Trunc_1|Trunc_2]). \\ & \underline{bitstring}([0,0]). & \underline{bitstring}([1,0,0]). & \underline{bitstring}([1,0,0,Trunc_1|Trunc_2]). \\ & \underline{bitstring}([0,1]). & \underline{bitstring}([1,0,1]). & \underline{bitstring}([1,1,0,Trunc_1|Trunc_2]). \\ & \underline{bitstring}([1,1]). & \underline{bitstring}([1,1,1]). & \underline{bitstring}([1,1,1,Trunc_1|Trunc_2]). \\ & \underline{bitstring}([1,1,1]). & \underline{bitstring}([1,1,1,Trunc_1|Trunc_2]). \\ &
```

In the case of bitstring example, no finite value D_{depth} that will cover the type without rounding up. While it is true that making D_{depth} larger produces a tighter analysis, it is unlikely to improve any downstream analyses. For example, the

 $[\]overline{^{131}}$ Note that this example does not require any input data.

space complexity analysis that we will see in §6.3.1 will still (correctly) predict that the amount of space required by fixpoint iteration is ∞ regardless of whether D_{depth} is 5 for 5 billion. ¹³²

We can obtain a more precise analysis that constrains the list argument to be well-formed. We use the following:

```
197 % Declare user parameters
198 params: b; bs.
199 % type suggestion for bitstring
200 bitstring(Xs) :- bs(Xs).
201 % propagation rules
202 b(0) <== true.
203 b(1) <== true.
204 bs([]) <== true.
205 b(X) <== bs([X|Xs]).
206 bs(Xs) <== bs([X|Xs]).</pre>
```

Notice that #200 initializes the type for bitstring(Xs). From there, abstract forward chaining proves that it is a valid upper bound. The output is

```
207 bit(0).
208 bit(1).
209 bitstring(Xs) :- bs(Xs).
```

It is no surprise that our symbolic evaluation method (abstract forward chaining) requires these kinds of termination strategies, as symbolic evaluation techniques, more broadly, are plagued by *path explosion*, and they employ truncation techniques (like ours) to make their methods practical (Baldoni et al., 2018, §5). State merging (Baldoni et al., 2018, §5) is another practical technique where two or more analyzer states are merged, i.e., replacing a single state representing their union, possibly a lossy union. Future work may wish to explore state-merging methods as a way to mitigate the aggressive truncation

¹³²In our experience as users of the type analyzer, we have seen that there are advantages to re-running the type analyzer with smaller values of the truncation parameters since they result in a smaller output that illustrates the same divergence issue.

that we do here. 133

In summary, we have described a truncation technique for ensuring that abstract forward chaining produces a valid upper bound in finite time. The bound is sound but may be quite loose. The hope is that a loose bound is more useful than a "type analysis did not converge" error. Users can inspect the system's output and provide additional constraint-propagation rules to avoid truncation. Thus, the type analysis system can be used as a theorem-proving assistant. ¹³⁴ Lastly, truncation allows downstream analyses to proceed, albeit with the loose type.

133 We gave a simple-type merging method in Algorithm 15 that may be useful for this purpose.

¹³⁴It would be interesting to explore the use of dynamic analysis as a means to inform static analysis. For example, a future system could synthesize plausible invariants that happen to be true of the program items when running on concrete data. From there, the system could recommend some assertions to the user that would tighten the static program analysis of the program.

6.2 Type-based Program Specialization

This section describes a system that specializes a given program to a specific input—output API. The system exploits the type analyzer of §6.1 to perform a program transformation (i.e., Dyna-to-Dyna code generation) that produces the specialized program. Consider the following example.

```
210 goal += a(I,I).

211 a(I,K) += b(I,J) * c(J,K).

212 goal += dead(X).

with the input type<sup>134</sup>

213 \overline{b}(-,-).

214 \overline{c}(-,-).

and output type<sup>134</sup>

215 \widehat{c}(-,-).
```

We can see that dead(X) is not defined; thus, the bottom rule (#212) provably never fires given the declared input type. Now, suppose we work backward from the outputs (i.e., goal in the example). In that case, we see that we only require a subset of a(I,K)—specifically, just its diagonal a(I,I). We call a(I,I) useful, whereas the off-diagonal elements are useless. Useful items should also be non-0. Usefulness analysis is a form of sensitivity analysis: an item is useful if it is non-0 and setting its value to 0 would change the value of any output item, and useless otherwise. Computing the precise set of useful items requires concrete input data, so much like §6.1, we will estimate a (symbolic)

¹³⁴Recall that the underscore (_) creates a fresh, unnamed variable.

¹³⁴Here the marking $\hat{\cdot}$ is like the $\bar{\cdot}$ marking from §6.1, except for *useful* types (see below).

¹³⁵In the parlance of finite-state machines and context-free grammars, a useful item is *both* accessible and co-accessible.

upper bound on the set of useful items for a specific input-output type.

Below, we exploit our knowledge of the useful items to manually *specialize* each rule of the example program so it computes goal more efficiently:

```
goal += a(I,I).

217 a(I,I) += b(I,J) * c(J,I). % only builds a's diagonal
```

Notice that this program is not fully equivalent to the previous one because the off-diagonal a(I,K) items are now 0-valued. However, the declared input-output API is preserved: the output values (namely, goal) are guaranteed to be equal to those of the original program.

We can improve this program further by *abbreviating* the a(I,I) items to a'(I) because they always carry a redundant variable:

```
218 goal += a'(I).
219 a'(I) += b(I,J) * c(J,I).
```

We can optionally include rules, such as a(I,I) += a'(I), that recover the original items from their abbreviated versions (assuming they are useful). Later, we will show examples of abbreviations that flatten nested terms, remove redundant variables, and drop variables that are never instantiated.

This section will develop an automated system for performing specialization transformations like these.

- First (§6.2.1), we describe extension to our type analysis method (§6.1) to estimate the set of useful items.
- Next (§6.2.2), we give a specialization method that transforms each rule—

¹³⁶In some parts of the thesis, this kind of semantics-altering transformation would require us to rename the specialized version a as it has a different meaning. To keep the exposition simple, we will allow it.

Algorithm 16 useful: For a given program \mathfrak{p} , input type $\overline{\mathfrak{D}}$, and output type $\widehat{\mathfrak{D}}$, define a new set of items $\widehat{\cdot}$ where \widehat{x} is true if x might be useful.

```
1. \triangleright Note: The function \widehat{\cdot} wraps a term in a novel function symbol.

2. \operatorname{def} useful (\mathfrak{p}, \overline{\mathfrak{D}}, \widehat{\mathfrak{D}}):

3. \widehat{\mathfrak{p}} \leftarrow \langle \mathcal{f} \rangle \triangleright New program

4. \widehat{\mathfrak{p}} \oplus = \overline{\mathfrak{p}} \oplus \overline{\mathfrak{D}} \triangleright Defines the set of the possible items \overline{\cdot}; same as §6.1

5. \widehat{\mathfrak{p}} \oplus = \widehat{\mathfrak{D}} \triangleright Base case: output types are useful

6. \operatorname{for}(x \oplus = y_1 \circ \cdots \circ y_K) \in \mathfrak{p}:

7. \operatorname{for} k \in [1:K]: \overline{137}

8. \widehat{\mathfrak{p}}[\widehat{y_k}] : -\widehat{x}, \overline{y_1}, \dots, \overline{y_K} \triangleright The subgoal y_k is useful if it appears in an active rule with a useful head
```

by expanding it into multiple versions per possible type—with the aim of reducing the total number of useless items and useless rule firings.

• Lastly (§6.2.3), we refine the specialization method to perform abbreviation, as in the example above.

6.2.1 Inferring a Type for Useful Items

To reason about the set of useful items for a specific API (i.e., a pair of an input type program $\widehat{\mathfrak{D}}$ and an output type program $\widehat{\mathfrak{D}}$, we transform the program \mathfrak{p} using Algorithm 16. The output will be a new program $\widehat{\mathfrak{p}}$ that defines new items marked with $\widehat{\cdot}$ such that \widehat{x} is true if x might be useful. Please review the annotated pseudocode below for a further explanation.

Performing type analysis on $\widehat{\mathfrak{p}}=\mathsf{useful}(\mathfrak{p},\overline{\mathfrak{D}},\widehat{\mathfrak{D}})$ using abstract forward chaining (§6.1) will infer a type for $\widehat{\mathsf{x}}$. Below is $\widehat{\mathfrak{p}}$

 $[\]overline{\,}^{137}$ It is safe to skip k when y_k is a constant or builtin.

```
% base case: output types are useful
                                                                            229 % booleanized program
                                                                            \overline{\text{goal}} := \overline{a}(I,I).
                                                                            \overline{goal} := \overline{dead}(X).
     % goal += a(I,I).
                                                                           \overline{a}(I,K) := \overline{b}(I,J), \overline{c}(J,K).
\widehat{a}(I,I) := goal, \overline{a}(I,I).
                                                                                  % input type
224 % goal += dead(X).
                                                                           \overline{b}(I,J).
225 \widehat{\text{dead}}(X) := \widehat{\text{goal}}, \overline{\text{dead}}(X).
                                                                           \overline{c}(J,I).
226 % a(I,K) += b(I,J) * c(J,K).
\widehat{b}(I,J) := \widehat{a}(I,K), \overline{b}(I,J), \overline{c}(J,K).
228 \hat{c}(J,K) := \hat{a}(I,K), \overline{b}(I,J), \overline{c}(J,K).
 The inferred type \widehat{\mathfrak{m}} = \operatorname{abstract} \operatorname{fc}(\widehat{\mathfrak{p}}) is
236 goal.
                         goal.
\widehat{a}(I,I).
                         \overline{a}(I,K).
238 b(I,J).
                         \overline{b}(I,J).
239 \hat{c}(J,K).
                         \overline{c}(J,K).
```

We see on #237 that only the diagonal a(I,I) is useful and the whole of a(I,K) is possible, as we observed in our manual analysis earlier.

Lastly, we note that $\widehat{\mathfrak{m}} \subseteq \overline{\mathfrak{m}}$ because we have added more constraints to the inference problem. We can see on line 16.8 that $\widehat{y_k}$ is only true if $\overline{y_k}$.

6.2.2 Specialization Transformation

Given a particular type $\widehat{\mathfrak{m}}$, ¹³⁸ we can specialize our program \mathfrak{p} by instantiating its rules against this type to see which rules will never (usefully) fire. Better yet, we can specialize each rule to its useful instantiations. We provide pseudocode in Algorithm 17.

Correctness.

 \hookrightarrow **Disjoint** $\widehat{\mathfrak{m}}$. We require the branches of $\widehat{\mathfrak{m}}$ to be pairwise disjoint. Disjointness guarantees that we can replace a rule r by nonoverlapping instanti-

¹³⁸ Note that the type $\overline{\mathbb{m}}$ from §6.1 would work too (swapping $\widehat{\cdot}$ with $\overline{\cdot}$ where it is needed); however, but the resulting specialized program would be less efficient since $\widehat{\mathbb{m}} \subseteq \overline{\mathbb{m}}$.

Algorithm 17 specialize: Specialize program \mathfrak{p} to the type $\widehat{\mathfrak{m}}$.

```
1. def specialize(\mathfrak{p},\widehat{\mathfrak{m}}):
2. assert the simple types in \widehat{\mathfrak{m}} are pairwise disjoint
3. \mathfrak{p}' \leftarrow \langle \zeta \rangle \triangleright New program
4. for (x \oplus = y_1 \circ \cdots \circ y_K) \in \mathfrak{p}:
5. for \langle \mathcal{C}_0, \mathcal{C}_1, \dots, \mathcal{C}_K \rangle \in \widehat{\mathfrak{m}}[\widehat{x}] \times \widehat{\mathfrak{m}}[\widehat{y}_1] \times \cdots \times \widehat{\mathfrak{m}}[\widehat{y}_K]: \triangleright Expand against the type
6. \mathcal{C} \leftarrow propagate(\mathcal{C}_0 \circ \mathcal{C}_1 \circ \cdots \circ \mathcal{C}_K)
7. if fail \in \mathcal{C}: continue \triangleright Safe to drop since it cannot fire
8. \mathfrak{p}'[x] \oplus = y_1 \circ \cdots \circ y_K \circ \mathcal{C} \triangleright The constraints \mathcal{C} become boolean side conditions
9. return \mathfrak{p}'
```

ations. This is needed to avoid double-counting. Given a disjoint useful type $\widehat{\mathfrak{m}}$, we can specialize the rules of the program and its different cases. The way that the specialization transformation works is to use the abstract instantiation of the rule to gather constraints of the types that it can take on. Using those types, we can add (concrete) constraints to the rule that partition it into its possible cases. If there are overlapping types, we can take an inferred type $\widehat{\mathfrak{m}}$ and make it disjoint by applying $\operatorname{disjoin}(\widehat{\mathfrak{m}})$ (Algorithm 9) as a preprocessing step.

 \hookrightarrow **User parameters required.** To leverage the user's parametric types, they must provide them in the input. Consider the two possible types form the $\gamma(X,Y)$ items in CKY (Example 18):

```
240 \overline{\gamma}(X,Y) := k(X), k(Y)
241 \overline{\gamma}(X,Y) := k(X), w(Y)
```

When the user provides a collection $\gamma(X,Y)$ we do not know which have k(Y) and which have w(Y). Thus, we require the user to provide these sets for our system to make their type distinctions.¹³⁹

We require the user to provide their parameters in addition to the input

¹³⁹In some cases, it may be possible to automatically determine the user's parameter sets from the input data. A warning is printed if the user forgets to provide at least one instance for each parameter.

data. For instance, the CKY example requires the user to provide methods that implement the n, k, and w constraints. This is done by defining them as items along with the input data. Since these items are constraints, they are $\{0,1\}$ -valued; only the 1-valued items need to be enumerated in the input data \mathfrak{D} . Note that the < constraint that appears in CKY's input type is a built-in relation provided by the system.

6.2.3 Abbreviation Transformation

We extend the specialization transformation (§6.2.2) to *abbreviate* the names of items as we did in the matrix-diagonal example earlier. Abbreviation is a scheme that partitions each item by its possible simple types. It defines new items that are named by their simple type and any (necessary) arguments. Pseudocode is provided in Algorithm 18.

Correctness. Abbreviation and specialization (§6.2.2) are correct under the same conditions.

Abbreviation vs. specialization. Unlike the specialization transformation, the abbreviation transformation can omit certain constraints from the abbreviated program because they are statically guaranteed to hold given the types in the body of the rule (see line 18.27). The specialization transformation

They can be defined using rules, e.g., if sentence positions (n) are encoded as the integers $\{1, \dots, 10\}$, we can write n(X) := 1 < X < 10.

 $^{\{1,\}ldots,10\}$, we can write $n(X):-1\leq X\leq 10$.

141 The order of the tuple $\langle X_1,\ldots,X_M\rangle$ is arbitrary, but it must be consistent for each t. Completely free variables may be omitted from the tuple.

¹⁴²When adding the rule on line 18.30, we must apply the multiplicity correction (§3.1.2.3) if any free variables are summed in the abbreviation. We compute the set of local free variables in the original rule to those in the abbreviated rule. If these sets of variables are unequal, we \sim multiply the body of the new rule by $(\infty \cdot \underline{1})$.

Algorithm 18 abbreviate: Abbreviate program \mathfrak{p} according to the type $\widehat{\mathfrak{m}}$.

```
1. def abbreviate(\mathfrak{p}, \widehat{\mathfrak{m}}):
         assert the simple types in \widehat{\mathfrak{m}} are pairwise disjoint
        \mathfrak{p}' \leftarrow ?
                                                                                                           ▶ Will be the output program
        \triangleright Map the item x to the abbreviated name for the type identified by t
         def abbrev<sub>t</sub>(\boldsymbol{x}):
            if x is a builtins and constants: return x
             (\mathbf{x}': -\mathcal{C}') \leftarrow \mathsf{fresh}(\widehat{\mathfrak{m}}_t)
             \langle \mathsf{X}_1, \dots, \mathsf{X}_M \rangle \leftarrow \text{the ordered}^{141} \text{ set of variables in } (x' : - \mathcal{C}')
             \theta \leftarrow \mathsf{unify}(x, x')
            return \theta(gen_t(X_1, ..., X_M))
                                                                               \triangleright Here gen<sub>t</sub> is a novel symbol for simple type t
        \triangleright API preservation rules
        for (x: -C) \in \widehat{\mathfrak{m}}:
             Assign a unique identifier t to the simple type (x : -C)
            if t is an input:
14.
                 ▷ Input items should not be abbreviated to preserve p's API. The rule below feeds
                     the input items to the abbreviated items for its type t.
                 \mathfrak{p}'[\mathsf{abbrev}_t(\boldsymbol{x})] \oplus = \boldsymbol{x} \circ \mathcal{C}
                                                       \sim Here, the constraints {\cal C} are required to ensure that the correct set
                                                          of input items are passed to this abbreviated item.
17
                 ▷ The rule below recovers the original item from the abbreviated item; Note: the
                    rule added below is useless in p' unless x is an output item.
                 \mathfrak{p}'[\boldsymbol{x}] \oplus = \mathsf{abbrev}_t(\boldsymbol{x})
19
        for (x \oplus = y_1 \circ \cdots \circ y_K) \in \mathfrak{p}:
            \textbf{for } \langle \mathcal{C}_0, \mathcal{C}_1, \dots, \mathcal{C}_K \rangle \in \widehat{\mathfrak{m}}[\widehat{\pmb{x}}] \times \widehat{\mathfrak{m}}[\widehat{\pmb{y}_1}] \times \dots \times \widehat{\mathfrak{m}}[\widehat{\pmb{y}_{\kappa}}] \colon
21
                 \mathcal{C} \leftarrow \text{propagate}(\mathcal{C}_1 \circ \cdots \circ \mathcal{C}_K)
                 if fail \in \mathcal{C}: continue
                 \mathcal{C}' \leftarrow \operatorname{propagate}(\mathcal{C}_0 \circ \mathcal{C})
                 if fail \in C': continue
                 \triangleright The head's type constraints must be checked, but the constraints in \mathcal C can be
                     removed, as they are guaranteed to be satisfied.
                 C'' \leftarrow C_0 \setminus C
                 \langle t_0, t_1, \dots, t_K \rangle \leftarrow \text{identifiers for the current simple types}
                 \triangleright Replace x and each y_k with its type-specific abbreviation. <sup>142</sup>
                 \mathfrak{p}'[\mathsf{abbrev}_{t_0}(\boldsymbol{x})] \oplus = \mathsf{abbrev}_{t_1}(y_1) \circ \cdots \circ \mathsf{abbrev}_{t_K}(y_K) \circ \mathcal{C}''
30
        return p'
31
```

could not do that because the subgoals are not split into their different types. So a subgoal with more than one simple type could satisfy the constraints in only *some* of the cases. This is not an issue in the case of abbreviation because the items are split into their disjoint simple types. Consider, again, CKY. The subgoal rewrite(X,Y) has two possible simple types: $\overline{\gamma}(X,Y) := k(X)$, k(Y) and $\overline{\gamma}(X,Y) := k(X)$, w(Y). So specialization cannot guarantee much about the constraint on Y, but abbreviation can because rewrite(X,Y) split into its two types. We provide an NLP example of abbreviate in Example 21.

Example 20. Consider the following program for the total weight of all paths between pairs of nodes in a weighted graph:

```
242 path(I,I).
243 path(I,K) += edge(I,J) * path(J,K).
```

Suppose the input and output types are

```
\overline{\mathfrak{D}} = \{ \text{params: n. } \overline{\text{edge}}(I,J) :- n(I), n(J). \} \ and \ \widehat{\mathfrak{D}} = \{ \text{path}(I,K) :- n(I), n(J). \}
```

Abbreviation produces the following program:

```
244 % map each input type to its abbreviation
245 edge<sub>0</sub>(I,J) += edge(I,J) * n(I) * n(J).
246 % recovery rules for useful output items
247 path(I,K) += path<sub>1</sub>(I,K).
248 % abbreviated rules
249 path<sub>1</sub>(I,I) += n(I).
250 path<sub>1</sub>(I,K) += edge<sub>0</sub>(I,J) * path<sub>1</sub>(J,K).
```

Notice that the nonground path(I,I) items are gone because the variable I is now constrained to n(I).

Now, suppose we introduce a type (a) for the useful starting states:

 $^{^{143}}$ That said, we at least know that k(X) holds in both cases so that we could drop that constraint. Our specialization transformation does not do this optimization, but a subsequent analysis could do it.

```
\overline{\mathfrak{D}} = \left\{ \mathsf{params: a; n. } \ \overline{\mathsf{edge}}(\mathtt{I},\mathtt{J}) \ :- \ \mathsf{n(I), n(J).} \right\} \ and \ \widehat{\mathfrak{D}} = \left\{ \mathsf{path(I,K) :- a(I), n(J).} \right\}
```

Under this API, abbreviate *produces*:

```
% map each input type to its abbreviation edge<sub>0</sub>(I,J) += edge(I,J) * n(I) * n(J).
% recovery rules for useful items path(I,K) += path<sub>1</sub>(I,K).
% abbreviated rules path<sub>1</sub>(I,I) += a(I) * n(I). path<sub>1</sub>(I,K) += edge<sub>0</sub>(I,J) * path<sub>1</sub>(J,K).
```

This version further constrains the path(I,K) items, as they must also satisfy a(I).

6.3 Automatic Space and Time Analysis

This section builds on the type analysis provided by the methodology of §6.1 to develop an automated method for analyzing the space and time complexity of a Dyna program $\mathfrak p$ given data $\mathfrak D$. Given the types $\overline{\mathfrak m}$ associated with $\mathfrak p$, the key is to find parametric upper bounds on their *cardinality*. We will then use these to predict a parametric $\mathcal O$ -bound on the running time of forward chaining execution of $\mathfrak p(\mathfrak D)$. The bounds of this section assume that the program is executed with the agenda-based fixpoint iteration algorithm of §5.

6.3.1 Space Analysis

To bound the space complexity of our fixpoint iteration algorithms for evaluating $\mathfrak{m}=\mathfrak{p}(\mathfrak{D})$, it suffices to bound the size of $\mathsf{nz}(\mathfrak{m})$ because these algorithms only have to store $\mathfrak{m}.^{144}$ We will build on §6.1's upper bound $\overline{\mathfrak{m}}\supseteq \mathsf{nz}(\mathfrak{m})$. We will use |s| to denote the true cardinality of a set s and $|s|\ge |s|$ to denote an upper bound; we extend the notation to programs by implicitly coercing them to sets through $\mathsf{nz}(\cdot)$. Our high-level approach is $|\overline{\mathfrak{m}}| \stackrel{\text{def}}{=} \sum_{s \in \overline{\mathfrak{m}}} |s| \ge \sum_{s \in \overline{\mathfrak{m}}} |s| \ge |\overline{\mathfrak{m}}| \ge |\mathfrak{m}|,^{145}$ where the size |s| of a simple type s counts the number of satisfying assignments

¹⁴⁴We assume that any additional data structures, such as those used to speed up the lookup operations, increase the space complexity by only a constant factor. Additionally, we assume that each value in $\mathbb V$ is constant-sized.

¹⁴⁵If simple types in $\overline{\mathfrak{m}}$ overlap considerably, then the last inequality will be loose because summation double-counts the overlapping regions. Removing subtypes (§6.1.2.4) addresses the most egregious case, but future work may consider transforming $\overline{\mathfrak{m}}$ into a *partition* of simple types as follows. Replace any overlapping pair of simple types s and t with $s \cap t$, $s \setminus t$, and $t \setminus s$. Repeat until no pairs s and t overlap. Set difference s can be implemented using negated versions of the elementary constraints—including equality (i.e., unification), which is usually represented implicitly.

to its variables, and $\lceil s \rceil$ is an upper bound on this.

It is well-known that counting the variable assignments that satisfy a set of constraints can be computationally expensive. However, the more important issue in our case is that we do not even know how the parametric constraints are defined. The remainder of this section considers how to construct a bound $\lceil s \rceil$ from some user-supplied information about the parametric constraints.

Input Size Specification. Consider the simple type $\overline{\gamma}(X,W) := k(X), w(W)$. A straightforward upper bound can be found:

$$|\bar{\gamma}(X,W) := k(X), w(W)| \leq [k(X)] \cdot [w(W)] = k \cdot w$$

where k and w are user-specified size parameters. (We will explain shortly how size parameters are specified.) However, matters become more complicated when there are multiple constraints on the same variable. Consider the simple type $\overline{\beta}(X,I,K) := k(X), \ n(I), \ n(K), \ I < K.$

It is clear that at most $k \cdot n \cdot n$ assignments can satisfy the first 3 constraints. Additionally, imposing the fourth constraint $\mathbb{I} < \mathbb{K}$ can only reduce that number—even though the fourth constraint in isolation would have infinitely many satisfying assignments ($|\mathbb{I} < \mathbb{K}| = \infty$).

We will describe how to derive an upper bound from the *sequence* of constraints c_1, \ldots, c_M in the body of the simple type s (including any constraints added by propagation, as they may improve the bound). Reordering the constraints does not affect the number of satisfying assignments; however, each of the M! orders may give a different upper bound, so we select the tightest of these.

The idea is to use **conditional cardinalities**. Let \mathcal{V} be any set of variables. While c may be satisfied by many assignments to its variables, only some are compatible with a *given* assignment to \mathcal{V} . Suppose the user supplies a number $\lceil c \mid \mathcal{V} \rceil$ such that no assignment to \mathcal{V} can be extended into more than $\lceil c \mid \mathcal{V} \rceil$ assignments to vars(c). The case $\lceil c \mid \mathcal{V} \rceil = 1$ is of particular interest: this is necessarily true if $\text{vars}(c) \subseteq \mathcal{V}$ and asserts a functional dependency otherwise.

For a simple type s of the form \bar{x} : $-c_1, \ldots, c_M$, we can obtain an upper bound on the cardinality of s: $|s| \leq \prod_{m=1}^{M} \lceil c_m \mid \mathcal{V}_{\leq m} \rceil$, where $\mathcal{V}_{\leq m} = \bigcup_{i=1}^{m-1} \operatorname{vars}(c_i)$. More generally, an upper bound on the *conditional* cardinality of s, given an assignment to some set of variables \mathcal{V} , is $\prod_{m=1}^{M} \lceil c_m \mid \mathcal{V} \cup \mathcal{V}_{\leq m} \rceil$. Minimizing this over all orderings gives our final upper bound $\lceil s \mid \mathcal{V} \rceil$ (with $\lceil s \rceil \stackrel{\text{def}}{=} \lceil s \mid \emptyset \rceil$).

But where do the bounds $\lceil c \mid \mathcal{V} \rceil$ come from? Consider the built-in constraint times, which is an infinitely large relation on three numbers. Because times has functional dependencies, we can take $\lceil \text{times}(X,Y,Z) \mid \mathcal{V} \rceil = 1$ whenever \mathcal{V} contains at least two of $\{X,Y,Z\}$. We allow a user to declare this conditional cardinality with the following notation:

```
259 |times(+X, +Y, Z)| \le 1.

260 |times(X, +Y, +Z)| \le 1.

261 |times(+X, Y, +Z)| \le 1.
```

For example, #261 declares that given values for X and Z, the times(X,Y,Z) constraint will be satisfied by at most one value of Y.

Given all such declarations, we obtain $\lceil c \mid \mathcal{V} \rceil$ as the min of the upper bounds provided by conditional cardinality declarations that unify with c and whose

Provided that all variables in the head \bar{x} are constrained somewhere in the body. If not, $|s| = \infty$.

variables marked with + are all in \mathcal{V} . This includes an implicit declaration that $\lceil c \mid \mathsf{vars}(c) \rceil \leq 1$. (Another implicit declaration is $\lceil \mathsf{fail} \rceil \leq 0$.)

The same notation is useful for bounding input relations, such as the word relation in Example 6.

```
\frac{1}{262} | \overline{\text{word}}(W, +I, K) | \le 1.
\frac{1}{263} | \overline{\text{word}}(W, I, +K) | \le 1.
\frac{1}{264} | \overline{\text{word}}(W, I, K) | \le n.
```

The first declaration says that there is at most one word w, paired with a single end position K, starting at each position K. The second declaration is the mirror image of this. The third says that there are at most K words overall in the input sentence—where the symbol K is a **size parameter**.

We may also use size parameters to bound the cardinality of type parameters:

```
|\mathsf{k}(\mathsf{X})| \leq k. |\mathsf{n}(\mathsf{I})| \leq n. |\mathsf{w}(\mathsf{W})| \leq w.
```

which bounds the number of nonterminals by k, sentence positions by n and terminals by w. Here k, n, and w are interpreted as symbols. We will combine these symbols to form a **parametric expression** that upper bounds the cardinality of the items in the program. For example, as we saw earlier, $\lceil \overline{\beta}(X,I,K) \rceil = kn^2$. After adding up the sizes of simple types, our system can compute an asymptotic upper bound on the resulting expression, e.g., $kn^2 + kn + k^2n = \mathcal{O}(kn \max(n,k))$.

A final trick is to specify the cardinality or conditional cardinality of a composite constraint, which can be matched against a subsequence of c_1, \ldots, c_M when constructing an upper bound. For example, $|even(X), prime(X)| \le 1$ says that there is only one even prime, while $|\overline{\gamma}(X,Y), k(Y)| \le g$ says that no nonterminal symbol X can rewrite as more than g different nonterminal symbols Y. (It does

not, however, constrain how many ways x can rewrite as a terminal symbol y because these do not satisfy the k(y) constraint.)

Limitations.

- Certain infinite relations, such as those that can be expressed with free variables and delayed constraints, have an estimated size of ∞ under this approach (cf. footnote 146). However, certain infinite relations can be represented efficiently by the solvers in §5 as finite-size programs. However, because our approach is based on bounding the number of non-0 items rather than the size of their representation, our size estimates end up extremely loose.
- We assume that all built-in relations ground out all variables that appear
 in them. Future work should integrate a more sophisticated mode system
 (e.g., Overton (2003) and Filardo (2017)) for reasoning about the instantiation
 states of variables. Mode analysis can also be used to reason about the set of
 delayed constraints, so it should help address the limitation mentioned earlier.
- Our space bounds are always (piecewise) polynomials in the input-size parameters; they cannot express exponential size bound, for example.

6.3.2 Time Analysis

Suppose that the Dyna program is executed by the agenda-based algorithm (Algorithm 8). This algorithm maintains two programs: the **chart** m and the **agenda** \eth . The chart represents the current estimate of the fixpoint (§4.3). The

agenda is a queue of updates to be applied to the values in \mathfrak{m} . When an update to y' is popped, its value is updated in the chart. This change is then propagated to other items x. The updates to these x items are placed on the agenda and carried out only later. Crucially, the affected x items are found by matching y' (the **driver**) to any subgoal of any rule $x \oplus = \cdots$, and then querying the chart for items (**passengers**) that match that rule's other subgoals.

Propagating an update from a given driver y' through a rule r requires us to find all the satisfying assignments to $\mathsf{body}(r)$ that are consistent with binding the driver subgoal to y'. Although r is not a type, a subgoal in its body such as $\gamma(X,Y)$ does impose a constraint on X and Y—namely the constraint $\overline{\gamma}(X,Y)$ from our type analysis (§6.1), which is true only if $\gamma(X,Y)$ could have a value in the chart.

Suppose we have matched a driver y' against a subgoal of r and have already matched 0 or more other subgoals to passengers retrieved from the chart. These matches have instantiated some variables $\mathcal V$ in the rule. For each subgoal match, we also know from our static type analysis (§6.1.2) that the variables involved in that match must have satisfied certain constraints $\mathcal C$. If $\gamma(\mathsf X,\mathsf Y)$ is the next subgoal that we will match against the chart, we will get at most $N \stackrel{\mathrm{def}}{=} [\mathcal C, \overline{\gamma}(\mathsf X,\mathsf Y) \mid \mathcal V]$ answers. Our methods from the previous section can compute this upper bound even though $\overline{\gamma}$ is not a built-in or parametric constraint: simply refine the item type $\overline{\mathfrak m}$ by unifying the heads of all simple types $s \in \overline{\mathfrak m}$ with $\overline{\gamma}(\mathsf X,\mathsf Y)$, and then upper-bound the conditional cardinality of this refined type $\overline{\mathfrak m}'$.

Thanks to the indexing of the chart, the runtime of getting and processing N answers is proportional to 1+N in the worst case (with the 1 representing the

overhead of attempting the query even if there are no answers). For each answer, we must then proceed to match the remaining passengers in all possible ways.

The pseudocode in Algorithm 19 makes this precise. Two subtleties are worth pointing out. First, $\overline{\mathbf{m}}$ is disjunctive. The constraints \mathcal{C} that we accumulate from previously matched passengers depend on the particular simple types that cover those passengers, so we must iterate through all possibilities. Second, while the driver subgoal is chosen by the agenda, the order in which to instantiate the remaining subgoals (passengers) is up to us. The pseudocode always chooses the next subgoal that will minimize the upper bound on runtime, given the simple types of the previously matched passengers. In other words, instead of following the fixed subgoal order stated in rule r, we are using our static analysis to optimize the subgoal ordering to have better running time. 147

Correctness and Limitations. The running time bound produced by Algorithm 19 assumes that the number of times an item drives is *constant*. If the program is acyclic, then the topologically ordered agenda (§5.1.4) guarantees that each item is popped at most once, so the bound is correct. However, if the program is cyclic, items may drive more than a constant number of times. In the cyclic case, we can interpret our bound as the worst-case complexity of one iteration of

¹⁴⁷Future work may wish to explore methods that directly optimize the average execution time of the passenger query order in place of our upper bound. The average execution time could be estimated and optimized using machine learning techniques. Vieira, Francis-Landau, et al. (2017) describe a general framework for such optimization. In their terminology, the specific passenger ordering would be a *strategy* for the *task* of answering the passenger query. Their method would additionally allow the ordering choice to be made dynamically by a *policy* that has access to the specific passenger query (i.e., it is not limited to just its type) and the status of the available indexing data structures.

Algorithm 19 runtime: Infer a runtime bound for p given a derived type $\overline{\mathfrak{m}}$.

```
1. def runtime(\mathfrak{p}, \overline{\mathfrak{m}}):
        t \leftarrow 0
                                                                                                                            \triangleright Accumulator
        for (\overline{y'}: -\mathcal{C}') \in \overline{\mathfrak{m}}:
                                                                                                                ▷ Possible driver types
           egin{aligned} \mathbf{for} & (x \oplus = y_1 \circ \cdots \circ y_K) \in \mathfrak{p} \colon \\ \mathbf{for} & k \in [1{:}K] \colon \end{aligned} & Matching rules against driver <math>y'
                    for \subseteq unify(y_{k}, y'):
                        ⊳ Add the cost of driving the current rule
                        t += \text{runtime'}(\underbrace{\{y_1, \dots, y_K\} - \{y_k\}}_{Pascengare}, \mathcal{C}', \text{vars}(y_k)) \triangleright Note: y' = y_k \ due \ to \ unify
        return t
10. > This recursive function (which can be memoized for efficiency) tells us the op-
       timal runtime for a set of subgoals (P), with C being the set of conjunctive con-
       straints that have been accumulated and V the set of bound variables.
11. def runtime(\mathcal{P}, \mathcal{C}, \mathcal{V}):
        if \mathcal{P} = \emptyset: return 1
                                                                                           ▷ Base case: no subgoals remaining
        best \leftarrow \infty
13
        > Optimize over the choice of next subgoal p from the set of remaining subgoals
        for p \in \mathcal{P}:
           t \leftarrow 0
                                                                     ▷ Accumulate total suffix runtime (given choice p)
            for \langle \mathcal{C}', \theta \rangle \in \overline{\mathfrak{m}}[\overline{p}]:
                                                                \triangleright Lookup simple types in \overline{\mathfrak{m}} that match p. Note: we
17
                                                                    require an explicit result stream here (§3.1.2.1).
                C'' \leftarrow \operatorname{propagate}(\boldsymbol{\theta}(C \circ C'))
                q \leftarrow \lceil \mathcal{C}'', \overline{p} \mid \mathcal{V} \rceil
                                                                 ▷ Lookup conditional cardinality for passenger p
                if q = 0: continue
                                                                                           ⊳ For efficiency, avoid recursive call
                t += q \cdot \text{runtime'}(\boldsymbol{\theta}(\mathcal{P} - \{p\}), \mathcal{C}'', \mathcal{V} \cup \text{vars}(p))
                                                                                                                                    \triangleright Recurse
            best min=t
22
        return 1 + best
23
```

seminaive (Algorithm 7) as a function of the input-size parameters rather than the total complexity of execution. To obtain the total complexity, we must adjust the bound to account for the number of times each simple type may pop, which may depend on the input-size parameters and the prioritization strategy used. 148,149

Other limitations: Because our runtime bounds are built on the space complexity analysis from §6.3.1, we inherit its limitations. In particular, our bounds

¹⁴⁸Discussion of approximate fixpoints and convergence rates is provided in §4.4.

¹⁴⁹Recall (§5.4) that a cyclical program can be solved more efficiently by decomposing it into strongly connected components. Future systems for reasoning about the running time of solving cyclical programs should include a static SCC decomposition analysis.

are (piecewise) multivariate polynomials in the input-size parameters. This means, for example, that our system cannot infer an exponential runtime bound; it would instead predict an infinite running time, as it is impossible to upper bound an exponential with a polynomial. We also mention that our runtime analyzer produces an infinite upper bound whenever an infinite-cardinality delayed constraint might appear in the chart, even though our execution algorithms do not incur an infinite runtime cost when working with them. Lastly, our bounds assume that semiring operations run in constant time.

 \hookrightarrow *Ideas for future work*: Static analysis of cyclicity or, more generally, the strongly connected components of a program would clearly be useful for improving our running time analysis. It would also be useful to have runtime bounds for the other algorithms in §5 (e.g., newton). Our general approach of simulating the work done in each iteration of these algorithms should translate straightforwardly. It would also be interesting to provide runtime bounds that are specialized to specific families of semirings (e.g., commutative idempotent semirings), as they may have finite-time convergence guarantees.

6.4 Applications

This section provides additional examples of our static analyzer.

Example 21. Kuhlmann, Gómez-Rodríguez, et al. (2011) introduce an exact dynamic program for arc-eager parsing (Nivre, 2003) when the scoring function is arc-factored (Eisner, 1996). A Dyna encoding of their algorithm is given below:

```
266 \beta(I,I^0,K) += word(I,K).

267 \beta(I,I^1,K) += word(I,K).

268 \beta(I,H^0,K) += \beta(I,H^0,J) * score(H \leftarrow H') * \beta(J,H'^0,K).

269 \beta(I,H^0,K) += \beta(I,H^0,J) * score(H \rightarrow H') * \beta(J,H'^1,K).

270 goal += \beta(0,0^0,N) * len(N).

271 params: word; score; len.
```

Naïvely, one could think that the arc-eager parser runs in $\mathcal{O}(n^5)$. However, a tighter analysis reveals it actually runs in $\mathcal{O}(n^3)$. To derive this, consider the following type declaration:

The derived types are

```
277 \overline{\beta}(I,I^0,K) := n(I), n(K).

278 \overline{\beta}(I,I^1,K) := n(I), n(K).

279 \overline{goal}.
```

What is surprising about the derived types is that, although the rules suggest that H in $\beta(I,H^AB,K)$ may independently range over n, it turns out that it is always equal to I because only types that unify with it are $\overline{\beta}(I,I^AB,K)$ and $\overline{\beta}(I,I^AI,K)$. Now, we consider the running time with the following input-size declarations:

```
|n(N)| \leq n.
```

The inferred runtime bound is $\mathcal{O}(n^3)$.

Abbreviation (below) successfully drops the unnecessary variable, which improves the running time by a constant factor. The program below is isomorphic to the one that Kuhlmann, Gómez-Rodríguez, et al. (2011) discovered manually. We have dropped the n(I) constraints to make the output more readable. We have chosen abbreviated names that are similar to their original ones; we use subscripts to indicate the simple type's identifier.

```
% Map each input type its abbreviation
    score_0(H,H') += score(H \rightarrow H').
                                                 these items are flattened (nested terms are remove)
    score_1(H,H') += score(H \leftarrow H').
    word_2(I,K) += word(I,K).
    len_3(N) += len(N).
285
286 % recovery rules
287 \beta(I,I^0,K) += \beta_4(I,K).
                                     these are useless recovery rules that can be dropped
    \beta(I,I^1,K) += \beta_5(I,K).
goal += goal<sub>6</sub>.
   % abbreviated program rules
    \beta_4(I,K) += word_2(I,K).
     \beta_4(H,K) += \beta_4(H,H') * score_0(H,H') * \beta_5(H',K).
    \beta_4(H,K) += \beta_4(H,H') * score_1(H,H') * \beta_4(H',K).
294 \beta_5(I,K) += word_2(I,K).
295 \beta_5(H,K) += \beta_5(H,H') * score_0(H,H') * \beta_5(H',K).
    \beta_5(H,K) += \beta_5(H,H') * score_1(H,H') * \beta_4(H',K).
    goal_6 += \beta_4(0,N) * len_3(N).
```

The inferred runtime bound is unchanged by the abbreviation transformation.

Example 22. Below we give another version of CKY (Example 6) that encodes the input sentence as a list, e.g., sentence([papa, ate, the caviar]). This version takes the encodes the weighted production rules as $\gamma(X,Y,Z)$ and $\gamma(X,Y)$ like before, but it does not require the sentence to be encoded with the word(W,I,K) or length(N) items. As before, the item goal sums the total weight of all complete parses.

```
params: input; \gamma. output: goal.

buffer(Ws) += sentence(Ws). % initialize the buffer

buffer(Ws) += buffer([W|Ws]). % shift

\beta(X,[W|Ws],Ws) += \gamma(X,W) * buffer([W|Ws]).

\beta(X,I,K) += \gamma(X,Y,Z) * \beta(Y,I,J) * \beta(Z,J,K).

\beta(X,I,K) += \gamma(X,Y) * \beta(Y,I,K).

goal += \beta(s,Ws,[]) * sentence(Ws).
```

We specify the following input type:

```
params: w; ws; k.
\overline{\gamma}(X,Y,Z) := k(X), k(Y), k(Z).
\overline{\gamma}(X,Y) := k(X), k(Y).
\overline{\gamma}(X,Y) := k(X), k(Y).
\overline{\gamma}(X,W) := k(X), w(W).
sentence(Ws) :- ws(Ws).
```

Recall that k is the set of nonterminals and w is the set of words. The types for

the grammar rules are the same as the version of CKY that we analyzed in §6.1.

We define the input as type ws, which is a list of words. When we add the following propagation rules to encode that each of suffix of our list of words is also a list of words.

```
310 ws([]) <== true.
311 ws(Xs) <== ws([X|Xs]).</pre>
```

Lastly, we add the following propagation rules

```
w(X) \le w(X|Xs].

w(X) \le w(X|Xs].

v(X) \le w(X|Xs].
```

The first rule says that the first element of a word list is a word, and the second rule says that the set of nonterminals k is disjoint from the set of words k.

Abstract forward chaining derives the following type for the derived items

```
\begin{array}{ll} & \overline{\text{buffer}}(\text{Xs}) := \text{ws}(\text{Xs}).\\ & \overline{\beta}(\text{X},\text{I},\text{K}) := \text{k}(\text{X}), \text{ws}(\text{I}), \text{ws}(\text{K}).\\ & \overline{\text{goal}}. \end{array}
```

The runtime bound is the same as in ordinary CKY, $O(k^3 n^3)$, except that instead of encoding the input length as the number of sentence positions

```
|n(I)| \leq n
```

we encode it as the number of suffixes in the input sentence

```
|ws(Ws)| \leq n.
```

The abbreviation transformation, in this case, is not interesting; all it does is change the names of the existing items with no benefit.

Summary

This chapter developed several techniques for reasoning about a given program independently from specific input data. We described our system for deriving symbolic upper bounds on the program's set of non-Q-valued items as a function of input data characterized by unknown, symbolic user parameters. We built on that technique to estimate the set of *useful* for answering queries about a specific set of output items. We then used the estimated set of useful items to *specialize* the program so that it computes them more efficiently by omitting useless items. We then turned our attention to space and time complexity analysis. We describe an approach for bounding the size of the set of non-Q items as a function of user-specified size parameters. We related this bound to the space complexity of agenda-based fixpoint iteration. We then gave a technique for analyzing the runtime complexity of agenda-based fixpoint iteration as a function of the input-size parameters. We will revisit the complexity bounds and the program specialization techniques in our program optimizer (§8).

Chapter 7

Program Transformations

In this chapter, we will develop a collection of meaning-preserving program-toprogram transformations. These transformations will be used to systematically transform a given initial program into a faster program that computes the same values in §8.

Intuitively, program transformations may be regarded as *reusable tricks* for transforming one set of recursive definitions (a Dyna program) into a different (possibly faster-to-evaluate) set of recursive definitions (a different Dyna program). Since a Dyna program is essentially a system of equations that defines a valuation function (§ 4), we may regard program transformations much like the mathematical identities used to manipulate systems of mathematical equations (e.g., to solve them or to prove a new identity). Program transformations may be applied manually or automatically.

Program transformations—unlike execution algorithms (§5)—can operate

without concrete input data.¹⁵⁰ Thus, many transformations correspond to steps of abstract interpretation (Cousot and Cousot, 1977) in that they simulate steps taken by some concrete solver. Thus, we introduce **input and output declarations**: we assume that the initial program in our transformation sequence declares some items as **input** and/or **output** items. The rest are considered **intermediate items**. Below is an example of CKY with input–output declarations.

```
319 \beta(X,I,K) += \gamma(X,Y,Z) * \beta(Y,I,J) * \beta(Z,J,K).

320 \beta(X,I,K) += \gamma(X,Y) * \beta(Y,I,K).

321 \beta(X,I,K) += \gamma(X,Y) * word(Y,I,K).

322 goal += \beta(s,0,N) * len(N).

323 input: word(X,I,K); \gamma(X,Y,Z); \gamma(X,Y); len(N).

325 output: goal.
```

To be considered *correct*, a program transformation must preserve the mapping from the valuation of the input items to the original program's valuation of the output items (§4). However, a program transformation can introduce, destroy, or alter intermediate items.¹⁵¹

Outline. In what follows, we will give an informal overview of the speedups that are achievable through program transformations—providing a more formal treatment of each transformation in a later section.

7.1	Informal Introduction	144
7.2	Semantics-Preserving Transformations	152
7.3	Definition-Introduction Transformation	156
7.4	Fold and Unfold Transformations	160

¹⁵⁰Static analysis (§6.1) are another example.

¹⁵¹We will formalize notions of transformation correctness and program equivalence in §7.2.

7.5	Hoisting Transformations	•	•	•	•	•	•	•	•	•	•	•	•	•	 •	•	•	•	•	•	•	186
7.6	Discussion																					201

7.1 Informal Introduction

This section provides an overview and general flavor of the types of transformations we will develop in this chapter. We defer formal details of each transformation to later sections. We also note that the method of 6.3 can be used to *automatically* derive the big- \mathcal{O} bounds on space and time complexity in the examples below.

Define and fold. We start our discussion with non-recursive folding, as we expect it will be familiar to most readers.

Example 23. The following program computes the total weight of length-4 paths in a graph with edge weights w:

```
326 total += w(Y_1, Y_2) * w(Y_2, Y_3) * w(Y_3, Y_4) * w(Y_4, Y_5).
```

By exploiting the distributive property, we can rewrite this program into one that reuses partial sums of products for each path-suffix length:

```
327 total += rest<sub>1</sub>(Y<sub>1</sub>).

328 rest<sub>1</sub>(Y<sub>1</sub>) += w(Y<sub>1</sub>,Y<sub>2</sub>) * rest<sub>2</sub>(Y<sub>2</sub>).

329 rest<sub>2</sub>(Y<sub>2</sub>) += w(Y<sub>2</sub>,Y<sub>3</sub>) * rest<sub>3</sub>(Y<sub>3</sub>).

330 rest<sub>3</sub>(Y<sub>3</sub>) += w(Y<sub>3</sub>,Y<sub>4</sub>) * rest<sub>4</sub>(Y<sub>4</sub>).

331 rest<sub>4</sub>(Y<sub>4</sub>) += w(Y<sub>4</sub>,Y<sub>5</sub>).
```

This program computes the same value for total, but it is significantly more efficient to evaluate than the original program. Specifically, the original program runs in time $\mathcal{O}(N^5)$ where N is the number of nodes in the graph, and

the transformed program runs in time $\mathcal{O}(N^2)$. This is an example of the definition transformation, which proposes the new rest_i definitions, and the fold transformation, which grafts these definitions into the program.

Readers may recognize that the junction tree algorithm exploits the same transformation scheme (Lauritzen and Spiegelhalter, 1988; Shafer and Shenoy, 1990) for inference in probabilistic graphical models and join-order optimization in database query planning (Chekuri and Rajaraman, 1997). Gildea (2011) discusses many uses of the junction tree algorithm in weighted deduction systems used in NLP. However, these approaches do not transform *recursive* computations—they are targeted at a fixed set of variables rather than an unbounded set as in recursive computations. ¹⁵²

Example 24. Recall that CKY (Example 6) has the following expensive rule

```
332 \beta(X,I,K) += \gamma(X,Y,Z) * \beta(Y,I,J) * \beta(Z,J,K).
```

This rule can be sped up by summing over the variable Y separately from K:

```
333 \beta(X,I,K) += tmp(I,J,X,Z) * \beta(Z,J,K).
334 tmp(I,J,X,Z) += \gamma(X,Y,Z) * \beta(Y,I,J).
```

The transformed rules reduce the running time from $\mathcal{O}(K^3N^3)$ to $\mathcal{O}(K^2N^3+K^3N^2)$ where K is the number of nonterminals, and N is the sentence length. ¹⁵³

Unfold. The unfold transformation is essentially the inverse of the fold transformation and corresponds to inlining code. It takes as input a specific

¹⁵²We note that GMTK (Bilmes and Zweig, 2002) provided a Junction tree optimizer that could optimize a recursive graphical model, but to our knowledge, their system did not statically optimize across recursive levels, as we do in this chapter.

¹⁵³In NLP, this kind of transformation is often referred to as a *hook trick* (Eisner and Satta, 1999a; Huang, Zhang, et al., 2005; Gildea, 2011).

subgoal $y_k \in \text{body}(\mathbf{R})$ of some rule \mathbf{R} . The goal is to replace y_k within \mathbf{R} by its definition; however, its definition may span multiple rules \mathbf{R}' . Therefore, we will remove \mathbf{R} from the program and replace it with a specialized version of \mathbf{R} for each rule \mathbf{R}' whose head unifies with y_k . These rules \mathbf{R}' define y_k —except in the special case where y_k matches any input items, in which case we cannot unfold y_k because its complete definition is not available. Note that unfolding can make some rules obsolete. 154

Example 25. Suppose the user provided the inefficient program below, which could have been obtained by folding Example 23 in a suboptimal manner.

```
335 goal += tmp<sub>1</sub>(X_1, X_4, X_5).

336 tmp<sub>1</sub>(X_1, X_4, X_5) += tmp<sub>2</sub>(X_1, X_4) * w(X_4, X_5).

337 tmp<sub>2</sub>(X_1, X_4) += tmp<sub>4</sub>(X_1, X_2) * tmp<sub>3</sub>(X_2, X_4).

338 tmp<sub>3</sub>(X_2, X_4) += w(X_2, X_3) * w(X_3, X_4).

339 tmp<sub>4</sub>(X_1, X_2) += w(X_1, X_2).
```

While the above program is correct, its runtime is $\mathcal{O}(N^3)$, which is worse than the $\mathcal{O}(N^2)$ version we saw before. It has no variables that can be eliminated, so there are no fold actions or new definitions that can improve its runtime complexity. To improve it, we first have to undo the poor choices. If we unfold Example 25 multiple times, we can arrive at the version in Example 23. From there, we can refold to get the same efficient version.

An unfold transformation will usually increase or preserve the program's runtime degree, but there are exceptions:

Example 26.

¹⁵⁴Recall that the methods in §6.2 can identify and remove rules that cannot fire based on the declared inputs or are unused by any of the declared outputs.

```
a(I,K) += b(I,J) * c(J,K).

trace += a(L,L).
```

Unfolding a(L,L) in #341 decreases the program's runtime from $O(N^3)$ to $O(N^2)$ assuming N upper bounds cardinality of any of variables in the program.

```
trace += b(L,J) * c(J,L).

a(I,K) += b(I,J) * c(J,K).
```

Thus, this is an example of an immediately useful unfold. The program, in this case, computes the trace of a matrix product, and the role of unfold is to specialize the sub-program that computes the entire matrix product a(I,K) to the site where the product is used, which only seeks its diagonal a(L,L). Such optimizations are easy for programmers to miss. ¹⁵⁵

Transforming recursion. In the following example, we transform an inductive definition—effectively eliminating a variable that is propagated through an (unbounded) recursion. Such a rewrite is impossible with only the tools from probabilistic graphical models and database query planning.

Example 27. We begin with the following program, which computes the total weight (goal) of all paths in a graph that start in a set start and end in a set stop.

```
goal += start(I) * path(I,K) * stop(K).

path(I,I).

path(I,K) += path(I,J) * edge(J,K).
```

The running time of this program is $\mathcal{O}(N^3)$, where N is the number of nodes in the graph. We will derive a more efficient $\mathcal{O}(N^2)$ program by replacing the intermediate path(I,K) items, which compute the total weight of all paths from I to

¹⁵⁵Many numerical libraries (e.g., Theano Development Team (2016), Abadi et al. (2016), and Smith and Gray (2018)) capture optimizations like this one. However, like the case of graphical models, these systems are limited because they do not reason about recursively defined computations.

K for all pairs of nodes I and K, with startpath(K), which computes the total weight of paths from any node in start to a node K.

Suppose we generate a new set of items that are—at least initially—defined by the following rule, which is appended to our program:

```
startpath(K) += start(I) * path(I,K).

Next, we fold \#344 by \#347:

348 goal += startpath(K) * stop(K).
```

As in Example 23, by introducing an intermediate item—in this case, startpath—we have managed to reduce the running time of #344 as its replacements #348 and #347 run in $\mathcal{O}(N+N^2)$ instead of $\mathcal{O}(N^3)$. Unfortunately, startpath still depends on path, which is currently the running time bottleneck for the program. Thus, we will attempt to rewrite startpath in terms of itself rather than path. Our next step is to unfold the path subgoal in the startpath rule (#347):

```
startpath(K) += start(K).
startpath(K) += start(I) * path(I,J) * edge(J,K).
```

Next, we recognize that start(I) * path(I,J) on #350 can be replaced by startpath(J) by folding according to its earlier definition (#347). We now have

```
startpath(K) += start(K).
startpath(K) += startpath(J) * edge(J,K).
```

Finally, a useless-rule elimination procedure (§6.2) detects that the path rules are not reachable from goal, which is our only declared output. Our final program is

```
goal += startpath(K) * stop(K).
startpath(K) += start(K).
startpath(K) += startpath(J) * edge(J,K).
startpath(I,I).
startpath(I,I) * edge(J,K).
```

The running time of our improved program is $\mathcal{O}(N^2)$.

Notice that not all folds are safe. For example, if instead of unfolding #347, we immediately folded the definition of startpath with itself, we would have the following ill-founded definition: startpath(K) += startpath(K). In §7.4, we will describe sufficient conditions to reject these bad folds and ensure safe folding.

Notice that to replace path in this example, we needed to generate a new relation startpath, which could replace it. Although startpath was initially defined in terms of path—through a few transformations—it became self-sufficient. Because discovering new definitions—like startpath—often require "divine inspiration;" such definitions are affectionately called Eureka definitions (Pettorossi and Proietti, 1994). Fortunately, there are many common strategies for proposing plausible Eureka definitions, such as joining existing subgoals as we did. We discuss the strategies used by our system in §7.3.

The next example follows the same general recipe as Example 27, except here we synthesize an efficient dynamic programming algorithm from an inefficient program that explicitly sums over the infinitely many paths in a graph where each path is specified as a list of its visited nodes.

Example 28. The following program sums the total weight of all paths in a graph. It is written in a very direct style as a list of nodes visited along the path.

```
358 goal += start(X) * \beta([X|Xs]).
359 \beta([X_1, X_2|Xs]) += edge(X_1, X_2) * \beta([X_2|Xs]).
360 \beta([X]) += stop(X).
```

As it is written, this program is horribly inefficient as there might be infinitely

many paths, which translates into infinitely many non- $\underline{0}$ β items. To speed this program up, we introduce the following definition, which sums the tail of the list x while exposing the first node x in the path:

```
b(X) += \beta([X|Xs]).
```

Next, we unfold #361 and obtain the following two rules that replace it:

```
362 b(X_1) += edge(X_1, X_2) * \beta([X_2|X_S]).
363 b(X) += stop(X).
```

We subsequently fold #362 by #361:

```
b(X) += b(Y) * edge(X,Y).
b(X) += stop(X).
```

Now, b is recursively defined in terms of itself rather than in terms of its original specification in #361 that made use of the inefficient β items. Moreover, there are finitely many non- $\underline{0}$ b items since there are finitely many nodes in the graph. Lastly, we fold #358 by #361 to graft b into the rule for goal:

```
goal += start(X) * b(X).
```

The final program is given below; we note that the β items are available, but the rules are now useless for deriving the goal, so the solver does not need to run them.

```
367 goal += start(X) * b(X).

368 b(X) += b(Y) * edge(X,Y).

369 b(X) += stop(X).

370 \beta([X,Y|Xs]) += edge(X,Y) * \beta([Y|Xs]).

371 \beta([X]) += stop(X).
```

Speculation and left-corner transformations. The example below briefly shows a useful rearrangement of recursive equations through the speculation transformation (§7.5).

Example 29 (Unary Cycle Factoring; adapted from Eisner and Blatz (2007)).

Consider again the CKY program (Example 6):

```
372 \beta(X,I,K) += \gamma(X,Y,Z) * \beta(Y,I,J) * \beta(Z,J,K).

373 \beta(X,I,K) += \gamma(X,Y) * \beta(Y,I,K).

374 \beta(X,I,K) += \gamma(X,W) * word(W,I,K).

375 goal += \beta(s,0,N) * len(N).
```

This program has an inefficiency: it repeats the same unary-chain closure for all spans $\langle I, K \rangle$. Our use of speculation will refactor the program (see below) so that unary chain closures may be precomputed (independently of the sentence) and reused at each span. This implementation trick—attributed to Stolcke (1995)—is standard in probabilistic parsing.

```
376 u(X,X') += \gamma(X,Y) * u(Y,X').

377 u(X,X') += \gamma(X,X').

378 \beta(X,I,K) += u(X,X') * \beta'(X',I,K).

379 \beta(X,I,K) += \beta'(X,I,K).

380 \beta'(X,I,K) += \gamma(X,Y,Z) * \beta(Y,I,J) * \beta(Z,J,K).

381 \beta'(X,I,K) += \gamma(X,W) * word(W,I,K).

382 goal += len(N) * \beta(S,0,N).
```

(A nearly identical program optimization can be achieved with the closely related left-corner transformation (§7.5.2).)

Next steps. In this section, we have provided a high-level idea of the program transformations that we develop in this chapter. In what follows, we will formalize each of the transformations and give further examples of how they are used and when they are useful. But first, we will formalize program equivalence and transformation safety.

7.2 Semantics-Preserving Transformations

We say that a program-to-program transformation $\mathfrak{p}^{(s)} \to \mathfrak{p}^{(t)}$ is **semantics-preserving** (or **safe**, for short) if the programs $\mathfrak{p}^{(s)}$ and $\mathfrak{p}^{(t)}$ are **semantically equivalent**, i.e., they define precisely the same valuation mapping (§4).¹⁵⁶ As a shorthand, we write $\mathfrak{p}^{(s)} \equiv \mathfrak{p}^{(t)}$ to assert (or stipulate) that the two programs are semantically equivalent. In other words, $\mathfrak{p}^{(s)} \equiv \mathfrak{p}^{(t)} \iff \mathfrak{p}^{(s)}[x] = \mathfrak{p}^{(t)}[x]$ for all $x \in \mathbb{H}$. We will also make use of the following stronger notion of semantic equivalence, which implies the former (Proposition 5).

Definition 31. We say that $\mathfrak{p}^{(s)}$ and $\mathfrak{p}^{(t)}$ are **strongly equivalent** if there exists a value-preserving bijection $\phi_{s \to t} : \mathfrak{p}^{(s)}\{\cdot\} \to \mathfrak{p}^{(t)}\{\cdot\}$ where $\mathfrak{p}^{(s)}\{\cdot\}$ and $\mathfrak{p}^{(t)}\{\cdot\}$ are the set of all derivations in $\mathfrak{p}^{(s)}$ and $\mathfrak{p}^{(t)}$, resepectively.

Definition 32. The function $\phi: \mathfrak{p}^{(s)}\{\cdot\} \to \mathfrak{p}^{(t)}\{\cdot\}$ is a **value-preserving bijection** if and only if for all $\delta' \in \mathfrak{p}^{(t)}\{\cdot\}$, there exists a <u>unique</u> $\delta \in \mathfrak{p}^{(s)}\{\cdot\}$ such that $\phi(\delta) = \delta'$, $[\![\delta]\!] = [\![\delta']\!]$, and that δ and δ' have the same root label.

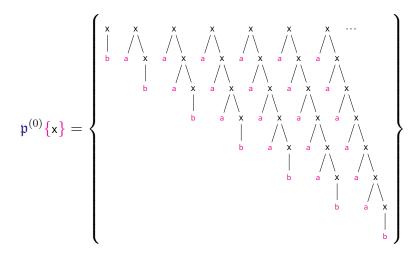
Intuitively, for each item x, the function $\phi_{s\to t}$ describes the translation from derivations $\mathfrak{p}^{(s)}\{x\}$ that use the rules in $\mathfrak{p}^{(s)}$ to derivations $\mathfrak{p}^{(t)}\{x\}$ that use the rules in $\mathfrak{p}^{(t)}$. However, value preservation restricts the translation such that the value contributed by each derivation of x is the same. Bijectivity ensures that no derivations are lost and no spurious derivations are created.

¹⁵⁶ Note that §4 does not directly discuss the semantics of partially specified programs, i.e., those with input relations. Fortunately, input relations can be reduced to built-in relations because they are simply items whose meanings are specified outside of the current program.

Example 30. Consider the geometric series program. We will illustrate how systematically different sets of derivations of different programs can index a common bag of contributions to the value of a given item x.

$$\mathfrak{p}^{(0)} = \begin{cases} 383 & x \neq b. \\ 384 & x \neq a \times x. \end{cases}$$

The value x is equal to the sum of a bag of (infinitely many) values, $[x] = b + ab + a^2b + a^3b + a^4b + a^5b + a^6b + \cdots$. Below, we show 7 (of the infinitely many) derivations of the item x. We see that the values of each derivation (i.e., the product of its leaves) correspond to the sequence of summands above.



The program $\mathfrak{p}^{(1)}$ below (obtained by unfolding #384) is strongly equivalent to the program above, $\mathfrak{p}^{(1)} \stackrel{\text{strong}}{\equiv} \mathfrak{p}^{(0)}$.

$$\mathfrak{p}^{(1)} = \begin{cases} 385 & x \neq b. \\ 386 & x \neq a \neq b. \\ 387 & x \neq a \neq a \neq x. \end{cases}$$

To see why we inspect the derivations of the item x. Below, we have ordered the derivations $\mathfrak{p}^{(1)}\{x\}$ so that the left-to-right ordering shows the correspondence (i.e., bijection) to $\mathfrak{p}^{(0)}\{x\}$.

By inspecting the leaves of the trees, we see that each of the values of the derivations can be aligned to exactly one of the original derivations (and vice versa). In other words, there is a value-preserving bijection between $\mathfrak{p}^{(0)}\{x\}$ and $\mathfrak{p}^{(1)}\{x\}$.

In the program below, we have systematically moved **b** up to the root of each derivation—leaving 1 in **b**'s original position.

$$\mathfrak{p}^{(2)} = \begin{cases} \frac{388}{389} & \text{x' += a * x'} \\ \frac{390}{390} & \text{x' += 1} \end{cases}.$$

$$\mathfrak{p}^{(2)} \{x\} = \begin{cases} \frac{x}{\sqrt{390}} & \frac{x}{\sqrt{390}}$$

Notice, that $\mathfrak{p}^{(2)}$ defines a new item x', which has the value $[x'] = 1 + a + a^2 + \cdots = a^*$. However, $\mathfrak{p}^{(0)}[x'] = 0$ because x' was not defined in $\mathfrak{p}^{(0)}$. To be semantics-preserving, we require that x' be defined in \mathfrak{p}_0 . Since the name x' is not taken in

 $\mathfrak{p}^{(0)}$ it is safe to add the following (unused) rules to \mathfrak{p}_0 that define the x' items.

$$\mathfrak{p}^{(0')} = \begin{cases} \frac{391}{392} & x \neq b. \\ \frac{392}{393} & x' \neq a \times x. \\ \frac{393}{394} & x' \neq \frac{1}{2}. \end{cases}$$

By adding the definitions to $\mathfrak{p}^{(0)} \not\equiv \mathfrak{p}^{(2)}$, we get $\mathfrak{p}^{(0')} \equiv \mathfrak{p}^{(2)}$. Augmenting the program with new definitions is what we call a definition introduction transformation (§7.3). Introducing new intermediate quantities is one of the most powerful methods at our disposal for devising a more efficient version of a given program. In general, whenever a transformation introduces a new definition, we ensure that the newly named items do not conflict with any other items in the universe of programs. Furthermore, we will add these definitions to all programs in the universe so they remain semantically equivalent.

Definition introduction. We treat definition introduction (§7.3) as in Tamaki and Sato (1984). When a transformation introduces a definition, we pretend it was introduced in the original program and propagate it throughout all programs derived from it. Under this scheme, the earlier programs would be retroactively augmented with the definition even though they did not technically have it when they were created. Maintaining a global symbol table for all programs in the universe is important to ensure that a transformation does not introduce a new definition that conflicts with any prior definition's name, including inputs and builtins.

Equivalence relation. As the name suggests, semantic equivalence (\equiv) is an equivalence relation. Thus, we can represent programs found through transformations using an undirected graph where the nodes are programs, and

an edge exists between programs $\mathfrak p$ and $\mathfrak p'$ if $(\mathfrak p \equiv \mathfrak p')$. However, since \equiv is reflexive, symmetric, and transitive, not all edges need to be explicit. Instead, $\mathfrak p$ and $\mathfrak p'$ are equivalent if an (undirected) path exists between them. Since testing the semantic equivalence of any pair of programs is undecidable, ¹⁵⁷ the graph of equivalent programs will be incomplete. In other words, we can only approximate it by including (directed) edges between programs that we can prove are equivalent (e.g., if a semantics-preserving transformation from $\mathfrak p$ to $\mathfrak p'$ exists, we add the edge)—this means that if an edge exists, we are certain the programs are equivalent, and if no edge, then the programs might be equivalent. Thus, any equivalence deduced by closure is sound.

7.3 Definition-Introduction Transformation

The definition-introduction transformation augments an existing program with rules defining some new set of items. We assume that the names of these items are distinct from the existing items. It is subsequently the job of the fold and unfold transformations to graft these new definitions into the program. We can often effectively replace the original definitions with new ones. Example 27 demonstrated how a new definition startpath could be transformed to be self-

¹⁵⁷Even testing the equivalence of a Dyna program *without variables* is undecidable. This is because Dyna programs can directly encode the language generated by a context-free grammar (CFG). Thus, determining if two Dyna programs are equivalent is at least as hard as deciding whether two context-free grammars generate the same language, which is an undecidable problem (Hopcroft et al., 2007).

sufficient and grafted into the main program by folding it into the goal rule—replacing the more expensive path relation. In general, discovering alternative definitions can require "divine inspiration." Thus, as we mentioned earlier, such rules are affectionately called *Eureka definitions* in many works in the program transformation literature (Pettorossi and Proietti, 1994).

The most common technique for proposing plausible Eureka definitions is the **tupling technique** (Pettorossi and Proietti, 1994).¹⁵⁸ Tupling is very general; it is parameterized as follows: Given a product of nonground terms $z_1 \circ \cdots \circ z_N$, and a set of variables $\{x_1, \ldots, x_L\}$ (which may be any subset of $\text{vars}(z_1 \circ \cdots \circ z_N)$ as well as fresh variables), the new definition formed by the tupling technique is then

$$\mathsf{tmp}(\mathsf{X}_1,\ldots,\mathsf{X}_L) \oplus = \boldsymbol{z}_1 \circ \cdots \circ \boldsymbol{z}_N$$

where tmp is a novel function symbol; the name tmp and the order of its arguments $x_1, ..., x_L$ are chosen arbitrarily. The new definitions proposed in the earlier examples of this chapter (23, 24, 27, and 28) were instances of a version of the tupling technique known as the variable-elimination technique (defined below).¹⁵⁹

The **variable-elimination technique** is an instance of the tupling technique where (1) we use existing subproducts in the program, and (2) we pick

¹⁵⁸Cohen, Simmons, et al. (2011) used a combination of tupling, fold, and unfold to give a general recipe for deriving programs that compute the weighted intersection of logic programs. Using this recipe, they derived several complex programs as the product of simpler programs (e.g., a CKY is the intersection of a WCFG and a WFSA).

¹⁵⁹This technique is sometimes called *loop absorption* (Pettorossi and Proietti, 1994). In NLP, this form of variable elimination is known as a *hook trick* (Eisner and Satta, 1999a; Huang, Zhang, et al., 2005; Gildea, 2011).

the minimal subset of variables from the chosen subproduct for which the new definition can correctly replace the subproduct. More formally, for a given program \mathfrak{p} , the variable-elimination technique takes a rule $R \in \mathfrak{p}$. Suppose that R is of the form $x \oplus = \alpha \circ \mu \circ \beta$ where α , μ , and β are sequence of subgoals. If \circ is commutative, then we may permute R's subgoals when matching this pattern. If any subset of μ 's variables does not appear in any of x, α , and β , then those variables may be locally summed in a new definition—providing a more efficient rule because it has fewer variables. The set of variables $\{x_1, \ldots, x_L\} = \text{vars}(\mu) \setminus (\text{vars}(\alpha) \cup \text{vars}(\beta) \cup \text{vars}(x))$ is the minimal set of variables for the new definition to fold into this position of \mathfrak{p} successfully. The remainder of the technique follows tupling.

We have already seen a number of examples of variable elimination in §7.1. In particular, Examples 23–25, 27 and 28 use variable elimination to propose new definitions that enable a speed-up. We highlight one of them below.

Example 31. We revisit Example 24 to see how it uses the variable elimination. Recall that Example 24 sped up the rule

```
395 \beta(X,I,K) += \gamma(X,Y,Z) * \beta(Y,I,J) * \beta(Z,J,K).
```

by transforming it into two rules where the variable \forall is summed separately from K:

```
396 \beta(X,I,K) += tmp(I,J,X,Z) * \beta(Z,J,K).
397 tmp(I,J,X,Z) += \gamma(X,Y,Z) * \beta(Y,I,J).
```

We will show how the rule #397 is the output of variable elimination. 160

 $[\]overline{^{160} \text{Below we padded the body with a } \underline{1}}$ to match the $(x \oplus \alpha \circ \mu \circ \beta)$ pattern.

$$\underbrace{\beta(X,I,K)}_{x} \; + = \underbrace{\frac{1}{\alpha}}_{\alpha} * \underbrace{\gamma(X,Y,Z) \; * \; \beta(Y,I,J)}_{\mu} * \underbrace{\beta(Z,J,K)}_{\beta}$$

We named our new relation tmp and determined its arguments as follows:

$$\{I,J,X,Z\} = vars(\gamma(X,Y,Z) * \beta(Y,I,J)) \setminus (vars(\underline{1}) \cup vars(\beta(Z,J,K)) \cup vars(\beta(X,I,K)))$$

We now see how the #397 is the output of the variable elimination technique. Note that we have to subsequently fold #395 by our new rule to get #396 to complete the example. Each definition introduction should be followed by fold to graft the new definitions into the program.

Later, we will see that the speculation and left-corner transformations (§7.5) are also ways to generate new definitions based on an existing program.

7.4 Fold and Unfold Transformations

This section describes a pair of transformations called fold and unfold. This pair of transformations has a long history in unweighted logic programming¹⁶¹. They were generalized to the semiring-weighted logic programming (Eisner and Blatz, 2007). At a high level, these transformations are analogous to applying the semiring distributive property (Definition 17) in either the expanding (unfold) or factoring (fold) direction. The folding transform is essential for grafting new definitions into the existing rules of a program and rewriting those definitions to be efficient self-supporting recursive definitions (as we have seen in $\S7.1$). We start with a basic version of the transformations ($\S7.4.1$) and, then, develop a generalized version (§7.4.2). The generalized transformations are novel in the weighted setting. The generalized transformations themselves are largely the same as in the unweighted setting. However, we had to overhaul the correctness proofs in order to prove that these transformations preserve our sum-over-derivations semantics (§4)—the existing theory for the unweighted setting only considered the preservation of at least one derivation rather than all derivations (Definition 31). As noted in the unweighted setting, the generalized fold and unfold transformations are quite powerful, but they require careful technical considerations to avoid mistakes.

¹⁶¹E.g., Tamaki and Sato (1984), Kawamura and Kanamori (1990), Bossi et al. (1992), Gergatsoulis and Katzouraki (1994), Pettorossi and Proietti (1994), Roychoudhury et al. (1999).

7.4.1 Basic Fold and Unfold

Unfold. At a high level, what the unfolding transformation $\mathfrak{p}^{(s)} \xrightarrow{\mathsf{unfold}} \mathfrak{p}^{(t)}$ does is to replace a given rule $r \in \mathfrak{p}^{(s)}$ with versions of that rule \mathcal{S} that correspond to expanding a specific subgoal against the rules that define that subgoal's value. The result of the fold is a new program $\mathfrak{p}^{(t)} = \mathfrak{p}^{(s)} \setminus \{r\} \cup \mathcal{S}$.

Example 32. Below, we unfold the first subgoal (q) in #398:

$$\begin{cases} 398 & p += q * t. \\ 399 & & \\ 400 & q += r. \\ 401 & q += s. \end{cases} \xrightarrow{unfold(\#398,1)} \begin{cases} 402 & p += r * t. \\ 403 & p += s * t. \\ 404 & q += r. \\ 405 & q += s. \end{cases}$$

Because there are two rules #400 and #401 that define q, we get two variants (#402 and #403) of #398 in the new program. In this example r= #398 and S= {#402, #403} Below, we show the derivations for p and q.

$$\left\{
\begin{array}{c|cccc}
p & p & q & q \\
/ \setminus & / \setminus & | & | \\
q & t & q & t & r & s \\
| & | & | \\
r & r & s & t & r & s
\end{array} \right\} \xrightarrow{\text{unfold}(\#398,1)} \left\{
\begin{array}{c|cccc}
p & p & q & q \\
/ \setminus & / \setminus & | & | \\
r & t & s & t & r & s
\end{array} \right\}$$

Notice how #398 is replaced in the first derivation with #402 and how #398 is replaced in the first derivation with #403. The two derivations of q are unchanged because none of the rules that appear in them were altered.

Fold. Folding is the inverse of unfolding. At a high level, a folding transformation $\mathfrak{p}^{(s)} \xrightarrow{\mathsf{fold}} \mathfrak{p}^{(t)}$ replaces a collection of rules $\mathcal{S} \subseteq \mathfrak{p}^{(s)}$ with a single rule r. In other words, the result of the fold is a new program $\mathfrak{p}^{(t)} = \mathfrak{p}^{(s)} \setminus \mathcal{S} \cup \{r\}$.

Example 33. Continuing Example 32, we use fold to invert the unfold transforma-

tion from Example 32. We simultaneously fold #406 by #408, and #407 by #409.

$$\begin{cases} 106 & p \ += \ r \ * \ t \, . \\ 407 & p \ += \ s \ * \ t \, . \\ 408 & q \ += \ r \, . \\ 409 & q \ += \ s \, . \end{cases} \xrightarrow{\text{fold}} \begin{cases} 410 & p \ += \ q \ * \ t \, . \\ 411 \\ 412 & q \ += \ r \, . \\ 413 & q \ += \ s \, . \end{cases}$$

This replaces the set of rules $S = \{\#406, \#407\}$ by the single rule r = #410. To check that we folded correctly, verify that unfolding #410 gives the program on the left (this is known as the reversibility test).

Definition 33. The **basic unfold transformation** takes $\langle \mathfrak{p}^{(s)}, r, j \rangle$ where $\mathfrak{p}^{(s)}$ is a program, r is a rule $(r \in \mathfrak{p}^{(s)})$, and j is an index of r's body $(0 \le j \le |\mathsf{body}(r)|)$. And produces a program $\mathfrak{p}^{(t)}$, which is given by the following procedure:

Algorithm 20 unfold: Apply the basic unfold transformation (Definition 33).

```
1.  \mathbf{def} \ \mathsf{unfold}(\mathfrak{p}^{(s)}, r, j) \colon 
2.  \left( \mathbf{x} \oplus \mathbf{y}_{[:j)} \circ \mathbf{y}_{j} \circ \mathbf{y}_{(j:]} \right) \leftarrow r 
3.  \mathfrak{p}^{(t)} \leftarrow \mathfrak{p}^{(s)} \smallsetminus \{r\} \qquad \qquad \triangleright \ \mathit{New program with rule r removed} 
4.  \triangleright \ \mathit{Add new rules by expanding subgoal } \ y_{j} \ \mathit{in rule r with its defining rules in } \ \mathfrak{p}^{(s)} 
5.  \mathfrak{p}^{(t)}[\mathbf{x}] \oplus \mathbf{y}_{[:j)} \circ \mathfrak{p}^{(s)}[\mathbf{y}_{j}] \circ \mathbf{y}_{(j:]} 
6.  \mathbf{return} \ \mathfrak{p}^{(t)}
```

Let S denote the set of rules added by the transformation. There is a one-to-one alignment $\sigma_{t \to s}: S \to \mathfrak{p}^{(s)}$; each rule in S maps to a distinct rule in $\mathfrak{p}^{(s)}$ that gave rise to it. The inverse of this alignment is denoted $\sigma_{s,t}: \mathfrak{p}^{(s)} \to S$. $\mathfrak{S}^{(s)}$

Definition 34. The basic fold transformation takes $\langle \mathfrak{p}^{(s)}, \mathcal{S}, r \rangle$ where $\mathfrak{p}^{(s)}$ is a program, $\mathcal{S} \subseteq \mathfrak{p}^{(s)}$, and r is a new rule. And, produces a program $\mathfrak{p}^{(t)} =$

¹⁶²The following behaviors are inherited from $\mathfrak{p}[\cdot]$. If subgoal j is a constant or input, then the unfold is a no-op. If subgoal j is a built-in relation, then it may evaluate if it is sufficiently instantiated. Recall from §3.1.2.3 that expanding a rule can lead to summing over a free variable. That case is handled here thanks to the result-stream machinery developed in §3.1.2.1.

 $(\mathfrak{p}^{(s)} \setminus S) \cup \{r\}$. However, the fold is valid if and only if it passes the **reversibility test**: $\exists j : \mathfrak{p}^{(s)} = \mathsf{unfold}(\mathfrak{p}^{(t)}, r, j)$. ¹⁶³

Algorithm 21 fold: Apply the basic fold transformation (Definition 34).

```
\begin{aligned} & \mathbf{def} \ \mathsf{fold}(\mathfrak{p}^{(s)}, \mathcal{S}, r) \colon \\ & \mathfrak{p}^{(t)} \leftarrow (\mathfrak{p}^{(s)} \setminus \mathcal{S}) \cup \{r\} \\ & \mathbf{assert} \ \exists j : \mathfrak{p}^{(s)} = \mathsf{unfold}(\mathfrak{p}^{(t)}, r, j) \\ & \mathbf{return} \ \mathfrak{p}^{(t)} \end{aligned} \qquad \triangleright \textit{Reversibility test}
```

Notice that the new rule r must be carefully synthesized for this transformation method to be useful and pass the reversibility test. The reversibility test was originally studied by Maher (1987).

Theorem 4. Suppose $\mathfrak{p}^{(s)} \xrightarrow{\text{unfold}} \mathfrak{p}^{(t)}$ according to Definition 33, then $\mathfrak{p}^{(s)} \stackrel{\text{strong}}{\equiv} \mathfrak{p}^{(t)}$. **Theorem 5.** Suppose $\mathfrak{p}^{(s)} \xrightarrow{\text{fold}} \mathfrak{p}^{(t)}$ according to Definition 34, then $\mathfrak{p}^{(s)} \stackrel{\text{strong}}{\equiv} \mathfrak{p}^{(t)}$.

Proofs of correctness in the unweighted setting are available in many works (e.g., Maher (1987)). However, the weighted setting (first proposed in Eisner and Blatz (2007)) lacks a formal proof of correctness for fold and unfold; we fill this gap in the literature in §B.2. Furthermore, our theoretical results strengthen Eisner and Blatz (2007)'s (informal) correctness claims by demonstrating *strong* equivalence of fold and unfold, rather than (weak) equivalence.

7.4.2 Generalized Fold and Unfold

In this section, we generalize the basic unfold transformation (Definition 33) to allow the rules used to expand the subgoal to come from a (possibly different)

¹⁶³For this test, syntactic equality is sufficient.

¹⁶⁴For a given r, we can search over $S \subseteq \mathfrak{p}^{(s)}$ to see if any of those folds pass the reversibility test.

program $\mathfrak{p}^{(d)} \stackrel{\text{strong}}{\equiv} \mathfrak{p}^{(s)}$. ¹⁶⁵ We call $\mathfrak{p}^{(d)}$ the **auxiliary definitions**. As we shall see in Example 34, auxiliary definitions give us additional flexibility because we can fruitfully mix and match semantically equivalent versions of a given program.

Example 34. Consider unfolding the geometric series $\mathfrak{p}^{(0)}$ to get $\mathfrak{p}^{(1)}$:

$$\mathfrak{p}^{(0)} = \begin{cases} \begin{smallmatrix} 414 & \text{x += b.} \\ 415 & \text{x += x * a.} \end{smallmatrix} & \xrightarrow{\text{unfold} \left(\#415,1\right)} & \mathfrak{p}^{(1)} = \begin{cases} \begin{smallmatrix} 416 & \text{x += b.} \\ 417 & \text{x += b * a.} \\ 418 & \text{x += x * a * a.} \end{cases} \text{ two a's}$$

So far, so good. What if we unfold again? Frustratingly, unfold $(\mathfrak{p}^{(1)}, \#418, 1)$ gives

which has unrolled the recurrence 4 levels instead of 3 levels. 166 Using generalized unfold, unfold $(\mathfrak{p}^{(1)}, ^{\#}418, 1, \mathfrak{p}^{(0)})$ with $\mathfrak{p}^{(0)}$ as auxiliary definitions, gives

which unrolls the recurrence to 3 levels. 167

Generalized fold. We also consider a generalized fold transformation, which relaxes the reversibility test in basic fold (Definition 34) to allow unfolding against auxiliary definitions in its reversibility test. Specifically, the generalized fold transformation modifies the reversibility test in Definition 34 from $\exists j: \mathfrak{p}^{(s)} = \mathsf{unfold}(\mathfrak{p}^{(t)}, r, j, \mathfrak{p}^{(t)})$ to $\exists j: \mathfrak{p}^{(s)} = \mathsf{unfold}(\mathfrak{p}^{(t)}, r, j, \mathfrak{p}^{(d)})$. This added flexibility enhances

¹⁶⁵Future work may want to explore other families of equivalence beyond strongly equivalent as strong equivalence is clearly not *necessary* for correctness.

¹⁶⁶Unfolding the bottom rule once more would unroll the recurrence up to 8 levels.

¹⁶⁷Unfolding the bottom rule once more (against $\mathfrak{p}^{(0)}$) would unroll the recurrence up to 4 levels.

the power of our transformation system: Example 27 and Example 28 would not be possible with only basic folds. We revisit Example 27 below.

Example 35. Recall Example 27, which computes the total weight (goal) of all paths in a graph that start in a set start and end in a set stop. Suppose we start with $\mathfrak{p}^{(0)}$ and unfold it to $\mathfrak{p}^{(1)}$:

```
\mathfrak{p}^{(0)} = \begin{cases} \frac{428}{429} & \text{goal} \ += \ \text{startpath}(\texttt{K}) \ * \ \text{stop}(\texttt{K}). \\ \frac{429}{430} & \text{path}(\texttt{I},\texttt{I}) \ += \ 1. \\ \text{path}(\texttt{I},\texttt{K}) \ += \ \text{path}(\texttt{I},\texttt{J}) \ * \ \text{edge}(\texttt{J},\texttt{K}). \end{cases}
\xrightarrow{\text{unfold}(\#429,2)} \mathfrak{p}^{(1)} = \begin{cases} \frac{432}{430} & \text{goal} \ += \ \text{startpath}(\texttt{K}) \ * \ \text{stop}(\texttt{K}). \\ \frac{433}{430} & \text{startpath}(\texttt{K}) \ += \ \text{start}(\texttt{I}) \ * \ \text{path}(\texttt{I},\texttt{J}) \ * \ \text{edge}(\texttt{J},\texttt{K}). \\ \frac{434}{430} & \text{startpath}(\texttt{K}) \ += \ \text{start}(\texttt{I}) \ * \ \text{path}(\texttt{I},\texttt{J}) \ * \ \text{edge}(\texttt{J},\texttt{K}). \end{cases}
```

Next, we use generalized folding to rewrite startpath more efficiently in terms of itself rather than path. We recognize that $S = \{ \#434 \}$ can be replaced by

```
r = (\text{startpath}(K) += \text{startpath}(J) * \text{edge}(J,K))
```

This transformation folds #434 by #429. Since #429 is in $\mathfrak{p}^{(0)}$, this transformation is a generalized fold. Specifically, the transformation fold $_{\{\#434\}\to r}^0$ gives

```
\mathfrak{p}^{(2)} = \begin{cases} ^{437} & \text{goal } += \text{startpath(K)} \ * \ \text{stop(K)}. \\ ^{438} & \text{startpath(K)} \ += \ \text{startpath(J)} \ * \ \text{edge(J,K)}. \\ ^{439} & \text{startpath(K)} \ += \ \text{startpath(J)} \ * \ \text{edge(J,K)}. \\ ^{440} & \text{path(I,I)} \ += \ 1. \\ ^{441} & \text{path(I,K)} \ += \ \text{path(I,J)} \ * \ \text{edge(J,K)}. \end{cases}
```

We deactivate the path rules as they are no longer used by goal. The running time of $\mathfrak{p}^{(2)}$ is $\mathcal{O}(N^2)$ where N is the number of nodes in the graph, which is considerably better than the $\mathcal{O}(N^3)$ running time of the initial program $\mathfrak{p}^{(0)}$.

Warning! The generalized fold and unfold transformations require some technical restrictions to be semantics-preserving. The definitions above are only suffi-

cient to guarantee **partial correctness**—meaning that we *can* lose derivations, but we *cannot* gain spurious derivations. To give a quick example, consider Example 35. If instead of unfolding #429, we immediately folded #429 by itself, we would arrive at the ill-founded definition: startpath(K) += startpath(K). In the remainder of this section, we will develop sufficient conditions to ensure specific fold and unfold transformations are semantics-preserving.

7.4.2.1 Formal Definitions

Generalized Unfold. Below, we generalize the basic unfold transformation (Definition 33) to allow auxiliary definitions ($\mathfrak{p}^{(d)}$, below) to be used. ¹⁶⁸

Definition 35. The generalized unfold transformation takes $\langle \mathfrak{p}^{(s)}, r, j, \mathfrak{p}^{(d)} \rangle$ where $\mathfrak{p}^{(s)}$ is a program, r is a rule $(r \in \mathfrak{p}^{(s)})$, j is an index of r's body $(0 \le j \le |\mathsf{body}(r)|)$, and $\mathfrak{p}^{(d)}$ is a program $(\mathfrak{p}^{(d)} \stackrel{\text{strong}}{=} \mathfrak{p}^{(s)})$. And produces a program $\mathfrak{p}^{(t)}$, which is given by the following procedure.

Algorithm 22 unfold: Apply the generalized unfold transformation (Definition 35).

```
1. \triangleright The pseudocode is similar to basic unfold (Algorithm 20). The difference is that we have introduced an additional argument \mathfrak{p}^{(d)}, and we use \mathfrak{p}^{(d)}[y_j] instead of \mathfrak{p}^{(s)}[y_j] on line 22.6.

2. def unfold(\mathfrak{p}^{(s)}, r, j, \mathfrak{p}^{(d)}):
```

4. $\mathfrak{p}^{(t)} \leftarrow \mathfrak{p}^{(s)} \setminus \{r\}$ \triangleright New program with rule r removed

5. \triangleright Add new rules where subgoal y_j in rule r with its defining rules in $\mathfrak{p}^{(d)}$

6. $\mathbf{p}^{(t)}[\mathbf{x}] \oplus \mathbf{y}_{[:j)} \circ \mathbf{p}^{(d)}[\mathbf{y}_j] \circ \mathbf{y}_{(j:j)}$

return $\mathfrak{p}^{(t)}$

^{3.} $\left(oldsymbol{x} \oplus oldsymbol{=} oldsymbol{y}_{[:j)} \circ oldsymbol{y}_{j} \circ oldsymbol{y}_{(j:]}
ight) \leftarrow r$

¹⁶⁸Pettorossi and Proietti (1994) have generalized unfold, but Roychoudhury et al. (1999) did not. Thus, it was missing in Roychoudhury et al. (1999)'s safety framework, which we have based ours on, and we will add it in §7.4.2.4.

Let S denote the set of rules the transformation adds. There is a one-to-one alignment $\sigma_{t \to d} : S \to \mathfrak{p}^{(s)}$; each rule in S maps to a distinct rule in $\mathfrak{p}^{(d)}$ that gave rise to it. The inverse of this alignment is $\sigma_{d \to t} : \mathfrak{p}^{(d)} \to S$.

Unlike basic unfolds (Definition 33), not all generalized unfolds are safe (cf. Example 36); sufficient conditions are presented in Theorems 1 and 2. Notice when $\mathfrak{p}^{(d)} = \mathfrak{p}^{(s)}$, we recover basic unfold. We often write $\mathfrak{p}^{(s)} \xrightarrow{\text{unfold} \frac{d,j}{r \to \mathcal{S}}} \mathfrak{p}^{(t)}$ to describe a specific unfolding relationship between $\mathfrak{p}^{(s)}$ and $\mathfrak{p}^{(t)}$. This provides a convenient shorthand to reference \mathcal{S} and r, which are functions of the tuple $\langle \mathfrak{p}^{(s)}, r, j, \mathfrak{p}^{(d)} \rangle$. When we omit the definition index d, e.g., $\mathfrak{p}^{(s)} \xrightarrow{\text{unfold} \frac{j}{r \to \mathcal{S}}} \mathfrak{p}^{(t)}$, we take d = s as a default value for the auxiliary definition argument.

Generalized Fold. Below, we generalize the basic fold transformation (Definition 34) to allow auxiliary definitions to be used in its reversibility test.

Definition 36. The generalized fold transformation works as follows: Given $\langle \mathfrak{p}^{(s)}, r, \mathcal{S}, \mathfrak{p}^{(d)} \rangle$ where r is a rule that replaces a subset of rules $\mathcal{S} \subseteq \mathfrak{p}^{(s)}$ and $\mathfrak{p}^{(d)} \stackrel{\text{strong}}{=} \mathfrak{p}^{(s)}$ are auxiliary definitions. The result is a program $\mathfrak{p}^{(t)} = (\mathfrak{p}^{(s)} \setminus \mathcal{S}) \cup \{r\}$. However, $\mathfrak{p}^{(t)}$ must satisfy the **generalized reversibility test**: $\exists j : \mathfrak{p}^{(s)} = \text{unfold}(\mathfrak{p}^{(t)}, r, j, \mathfrak{p}^{(d)})$.

```
Algorithm 23 fold: Apply the generalized fold transformation (Definition 36).
```

```
 \begin{aligned} \mathbf{def} & \text{ fold}(\mathfrak{p}^{(s)}, r, \mathcal{S}, \mathfrak{p}^{(d)}) : \\ & \mathfrak{p}^{(t)} \leftarrow (\mathfrak{p}^{(s)} \smallsetminus \mathcal{S}) \cup \{r\} \\ & \mathbf{assert} \ \exists j \colon \mathfrak{p}^{(s)} = \mathsf{unfold}(\mathfrak{p}^{(t)}, r, j, \mathfrak{p}^{(d)}) \\ & \mathbf{return} \ \mathfrak{p}^{(t)} \end{aligned} \qquad \qquad \triangleright \textit{Replace rules } \mathcal{S} \textit{ by the rule } r
```

Unlike basic folds (Definition 34), not all generalized folds are safe (cf. Example 25); sufficient conditions are presented in Theorems 1 and 3. Similar

to our notation for unfold, we often write $\mathfrak{p}^{(s)} \xrightarrow{\text{fold } d, j \\ S \to r} \mathfrak{p}^{(t)}$. Notice that taking $\mathfrak{p}^{(d)} = \mathfrak{p}^{(s)}$ does *not* recover the basic fold. Specifically, the generalized reversibility condition would not be the same because the basic fold would unfold with respect to $\mathfrak{p}^{(t)}$, and the generalized fold would unfold with respect to $\mathfrak{p}^{(d)}$, which equals $\mathfrak{p}^{(s)}$ rather than $\mathfrak{p}^{(t)}$.

7.4.2.2 Challenges

Unfortunately, not all generalized unfolds and folds are semantics-preserving.

We provide examples of each in this section.

Example 36 (Unsafe generalized unfold). We start with

$$\mathfrak{p}^{(s)} = \begin{cases} 442 & p \neq q. \\ 443 & q \neq r. \\ 444 & r \neq 1. \end{cases} \qquad \mathfrak{p}^{(d)} = \begin{cases} 445 & p \neq 1. \\ 446 & q \neq r. \\ 447 & r \neq p. \end{cases}$$

Clearly, $\mathfrak{p}^{(s)} \equiv \mathfrak{p}^{(t)}$ as they both have $[\![\mathfrak{p}]\!] = [\![\mathfrak{q}]\!] = [\![\mathfrak{q}]\!] = 1$. Furthermore, $\mathfrak{p}^{(s)} \stackrel{\text{strong}}{\equiv} \mathfrak{p}^{(t)}$ as we can see that there is a value-preserving bijection between

$$\mathfrak{p}^{(s)}\{\cdot\} = \left\{egin{array}{c|ccc} \mathfrak{p} & \mathfrak{q} & \mathfrak{r} \ dots & dots \ \mathfrak{q} & \mathfrak{r} & 1 \ dots & dots \ \mathfrak{p} & 1 \ dots & dots \ \mathfrak{p} & 1 \ dots & 1 \ dots \ \mathfrak{p} & 1 \ \ \ \mathfrak{p} & 1 \ \$$

However, unfolding #443 against $\mathfrak{p}^{(d)}$ gives the following incorrect program:

$$\mathfrak{p}^{(t)} = \begin{cases} 448 & \mathsf{p} \, + = \, \mathsf{q}. \\ 449 & \mathsf{q} \, + = \, \mathsf{p}. \\ 450 & \mathsf{r} \, + = \, 1. \end{cases} \qquad \mathfrak{p}^{(t)} \llbracket \cdot \rrbracket = \begin{cases} 451 & \mathsf{p} \, + = \, \emptyset. \\ 452 & \mathsf{q} \, + = \, \emptyset. \\ 453 & \mathsf{r} \, + = \, 1. \end{cases} \qquad \mathfrak{p}^{(t)} \{ \cdot \} = \left\{ \begin{array}{c} \mathsf{r} \\ | \\ | \\ 1 \end{array} \right\}$$

Inspecting the rules, we see an ill-founded cycle between #448 and #449 that causes its valuation to differ and derivation set to shrink in size.

Example 37 (Unsafe generalized fold). We start with ¹⁶⁹

$$\mathfrak{p}^{(0)} = \begin{cases} 454 & p \neq q \\ 455 & q \neq 1 \end{cases}$$

$$\mathfrak{p}^{(0)} \{\cdot\} = \begin{cases} q & p \\ | & | \\ 1 & q \\ | & | \end{cases}$$

Since [p] = [q] = 1, we might try to replace all instances of q on the right-hand side with p. So we define the following new program $\mathfrak{p}^{(*)}$ by folding with $S = \{\#454\}$, r = (p += p), j = 0, and definitions $\mathfrak{p}^{(d)} = \mathfrak{p}^{(0)}$. The output is the following (partially correct) program.

$$\mathfrak{p}^{(*)} = \begin{cases} 456 & p \neq p \\ 457 & q \neq 1 \end{cases}$$

$$\mathfrak{p}^{(*)}\{\cdot\} = \begin{cases} q \\ | \\ 1 \end{cases}$$

Although, the generalized reversibility condition is satisfied (unfold($\mathfrak{p}^{(0)}$, #456, 1, $\mathfrak{p}^{(0)}$) = $\mathfrak{p}^{(0)}$), this transformation has lost the derivation of \mathfrak{p} . At a high level, what has happened is that we have created an ill-founded definition for \mathfrak{p} ; this buggy definition is cyclical without any base cases.

7.4.2.3 What Causes Partial Correctness?

To better understand the fold and unfold transformations, we will construct their derivation mappings. Recall from Definition 31 that the derivation mapping is a function ϕ that maps derivations from $\mathfrak{p}^{(s)}\{\cdot\}$ into derivations from $\mathfrak{p}^{(t)}\{\cdot\}$. We require that ϕ is a value-preserving bijection from the transformation $\mathfrak{p}^{(s)} \longrightarrow \mathfrak{p}^{(t)}$ to be strongly semantics-preserving. What is interesting about the generalized fold/unfold case is the auxiliary definitions' role in construct- $\overline{}^{169}$ This is a Dyna adaptation of Example 4.1 from Pettorossi and Proietti (1994).

ing these mappings. Roughly speaking, the derivation mappings $\phi_{s \to t}$ will pivot through $\mathfrak{p}^{(d)}\{\cdot\}$ using $\phi_{s \to d}$ and $\phi_{s \to d}$. The challenge of pivoting is that it can result in infinite recursion, corresponding to (possibly) losing a derivation (i.e., partial correctness). We will use what we learn from the derivation mapping's failure modes to develop sufficient conditions for when the transform is safe (i.e., when the mapping is a value-preserving bijection).

We direct the reader to Algorithm 24 and Algorithm 25, which provide pseudocode for the derivation mappings of the generalized unfold and fold transformations. We wish to emphasize two things: (1) the derivation mappings are never run as part of the transformations—they are just used in our analysis of the transformations, and (2) neither Algorithm 24 or Algorithm 25 are guaranteed to define a value-preserving bijection—in fact, they are not even guaranteed to be functions, as they do not terminate for all inputs.

To see why the transformations go awry, we will study how the derivation mapping chokes on the following example:

Example 38. Continuing with the unsafe fold example (Example 37). Recall:

And we saw that the derivation $\delta^{(s)}$ is lost in this transformation. Consider the evaluation of $\phi_{s \to t}(\delta^{(s)})$ through Algorithm 25. Since the top rule of $\delta^{(s)}$ is

 $\mathbf{S}^{(\ell)} = (\mathsf{p} + \mathsf{q}) \ and \ \mathbf{S}^{(\ell)} \in \mathcal{S}, \ we \ fall \ in \ the \ bottom \ case \ (line \ 25.7).$

$$\phi \begin{pmatrix} \varphi \\ | q \\ | 1 \end{pmatrix} = \phi \begin{pmatrix} u^{-1} \begin{pmatrix} q \\ | 1 \end{pmatrix} \end{pmatrix}$$

$$\phi \begin{pmatrix} u^{-1} \begin{pmatrix} q \\ | 1 \end{pmatrix} \end{pmatrix}$$

The invocation of u^{-1} returns¹⁷⁰

$$\begin{pmatrix} p \\ | \\ q \\ | \\ 1 \end{pmatrix} = u^{-1} \begin{pmatrix} q \\ | \\ 1 \end{pmatrix}$$

Now, when we return to line 25.7, we observe the problem:

$$\phi \begin{pmatrix} p \\ q \\ l \\ 1 \end{pmatrix} = \begin{pmatrix} p \\ p \\ l \\ s \to t \end{pmatrix} = \begin{pmatrix} p \\ p \\ p \\ l \\ s \to t \end{pmatrix} = \begin{pmatrix} p \\ p \\ l \\ q \\ l \\ p \end{pmatrix}$$

$$\phi \begin{pmatrix} p \\ l \\ q \\ l \\ l \end{pmatrix} = p$$

$$\phi \begin{pmatrix} p \\ l \\ q \\ l \\ l \end{pmatrix} = p$$

$$\vdots$$

This recurrence never terminates as ϕ immediately calls itself with the same argument. It tries to build an infinitely tall tree and never succeeds because it does not hit base cases. Thus, the derivation $\delta^{(s)}$ is lost by the derivation mapping—meaning that our procedure for ϕ cannot find any corresponding $\delta^{(t)}$ in $\mathfrak{p}^{(t)}\{\cdot\}$. This means that we have failed to define a value-preserving bijection (ϕ is not even a function). Therefore, we cannot promise safety for this instance of generalized fold, nor should we!

To see why, we trace the steps of u^{-1} . Since $\mathfrak{p}^{(d)} = \mathfrak{p}^{(s)}$ the $\underset{s \to t}{\phi}$ (and $\underset{t \to s}{\phi}$) transformations are no-ops. The folder rule $\mathbf{D}^{(\ell)} = (\mathsf{p} + \mathsf{e} + \mathsf{q})$ (#454). So when we reconstruct, we get $y_j = \mathsf{p}$. We also have the following alignment: $\underset{d \to t}{\sigma} (\mathbf{D}^{(\ell)}) = (\mathsf{p} + \mathsf{e} + \mathsf{q})$.

The lesson learned from this example is that for the derivation maps (Algorithm 24 and Algorithm 25) to be correct, we need to ensure that cyclical recursion—like we just saw in Example 38—will not occur. We will use the nontermination of ϕ as a diagnostic for when generalized fold/unfold goes awry. When $\phi_{s\to t}(\delta)$ (Algorithm 24 and Algorithm 25) fails to terminate on some $\delta\in\mathfrak{p}^{(s)}\{\cdot\}$, the problematic δ corresponds to a derivation that might be lost. Thus, to ensure correctness, we need to reject transformations for which $\phi_{s\to t}$ might not always terminate. More formally, we have the following theorems, which are proven in §B.2:

Theorem 2. Given $\mathfrak{p}^{(s)} \xrightarrow{\text{unfold} \frac{d,j}{r \to S}} \mathfrak{p}^{(t)}$ and $\mathfrak{p}^{(s)} \stackrel{\text{strong}}{\equiv} \mathfrak{p}^{(d)}$. Suppose the function $\phi_{s \to t} \left(\delta^{(s)} \right)$, as specified in Algorithm 24, terminates for all $\delta^{(s)} \in \mathfrak{p}^{(s)} \{\cdot\}$. Then, $\phi_{s \to t}$ is a value-preserving bijection between $\mathfrak{p}^{(s)} \{\cdot\}$ and $\mathfrak{p}^{(t)} \{\cdot\}$. Furthermore, $\mathfrak{p}^{(s)} \stackrel{\text{strong}}{\equiv} \mathfrak{p}^{(t)}$.

Theorem 3. Given $\mathfrak{p}^{(s)} \xrightarrow{\text{fold} \stackrel{d,j}{S \to r}} \mathfrak{p}^{(t)}$ and $\mathfrak{p}^{(s)} \stackrel{\text{strong}}{\equiv} \mathfrak{p}^{(d)}$. Suppose the function $\phi_{s \to t} \left(\delta^{(s)} \right)$, as specified in Algorithm 25, terminates for all $\delta^{(s)} \in \mathfrak{p}^{(s)}\{\cdot\}$. Then, $\phi_{s \to t}$ is a value-preserving bijection between $\mathfrak{p}^{(s)}\{\cdot\}$ and $\mathfrak{p}^{(t)}\{\cdot\}$. Furthermore, $\mathfrak{p}^{(s)} \stackrel{\text{strong}}{\equiv} \mathfrak{p}^{(t)}$.

The termination condition manifests in our proof of correctness (§B.2); termination ensures that a well-founded ordering exists for our inductive proofs that establish value-preserving bijectivity. Verifying that these derivation mappings always terminate can be difficult. We describe a relatively simple system in §7.4.2.4 for ensuring termination that is based on static analysis.

Algorithm 24 Derivation mapping for $\mathfrak{p}^{(s)} \xrightarrow{\mathsf{unfold}_{r \to \mathcal{S}}^{d,j}} \mathfrak{p}^{(t)}$.

- 1. > The pseudocode assumes that rules are distinct. This can be arranged by marking each rule with an identifier that is checked by equality tests. When a rule is copied from one program to another, its identifier is unchanged. When a rule is created, it is assigned a unique identifier.
- 2. \triangleright We use derivation mappings ϕ and ϕ that establish $\mathfrak{p}^{(s)} \stackrel{\mathrm{strong}}{\equiv} \mathfrak{p}^{(d)}$. They must be fixed for all $\delta^{(s)}$.
- 3. \triangleright In the code below, α , μ , and β are sequences.
- 4. $\operatorname{def} \phi_{s \to t} \left(\boldsymbol{\delta}^{(s)} \right)$:
- case constant or builtin: **return** $\delta^{(s)}$

6. **case**
$$\bigwedge^{y_1} \circ \cdots \circ \bigwedge^{y_K}$$
: **return** $\phi \left(\bigwedge^{y_1} \right) \circ \cdots \circ \phi \left(\bigwedge^{y_K} \right)$

7. **case**
$$\begin{pmatrix} x \\ | \\ \mu \\ \triangle \end{pmatrix}$$
 and its top rule $\in \mathfrak{p}^{(t)}$: **return** $\phi \in \mathfrak{p}^{(t)}$

8. **case**
$$\begin{pmatrix} x \\ \alpha \\ y_j \\ \beta \end{pmatrix}$$
 and its top rule $\notin \mathfrak{p}^{(t)}$:

9. return

$$\phi \left(\bigwedge^{\alpha} \right) \quad \phi \left(u \begin{pmatrix} y_j \\ \triangle \end{pmatrix} \right) \quad \phi \left(\bigwedge^{\beta} \right) \\
s \to t \begin{pmatrix} \beta \\ \triangle \end{pmatrix}$$

- ightharpoonup The top rule of the transformed derivation is given by the 2^{nd} element of the pair returned from u (notation abuse). The call to u replaces the y_j subtree with a value-equivalent derivation in $\mathfrak{p}^{(s)}\{\cdot\}$ (the 1^{st} element of the pair), which is then recursively transformed.
- 10. \triangleright Pivoting function: Returns pair of a value-equivalent derivation and the new top-level rule in $\mathfrak{p}^{(t)}$. See Lemma 2 for a formal statement and proof.
- 11. **def** $u\left(\bigwedge^{y_j} \right)$:

 \triangleright Note: This function references $\langle s, t, d \rangle$ from the calling context.

12.
$$\boldsymbol{\delta}_{1}^{(s)} \leftarrow \overset{\boldsymbol{y}_{j}}{\wedge}; \quad \boldsymbol{\delta}_{1}^{(d)} \leftarrow \overset{\boldsymbol{\phi}}{\wedge} \left(\boldsymbol{\delta}_{1}^{(s)}\right)$$

 \triangleright Convert subtree into $\mathfrak{p}^{(d)}\{\cdot\}$.

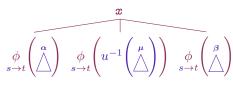
 \triangleright Remove the intermediate item y_i

- 14. $\mathbf{D}^{(\ell)} \leftarrow \text{top rule of } \boldsymbol{\delta}_1^{(d)}$
- 15. $\mathbf{T}^{(\ell)} \leftarrow \sigma(\mathbf{D}^{(\ell)})$
- 6. $\boldsymbol{\delta}_{2}^{(s)} \leftarrow \overset{a \to t}{\phi} \left(\boldsymbol{\delta}_{2}^{(d)} \right)$
- 17. **return** $\left< oldsymbol{\delta}_2^{(s)}, \mathbf{T}^{(\ell)} \right>$

- ▷ Use the alignment to retrieve the new top rule
 - \triangleright Convert back into $\mathfrak{p}^{(s)}\{\cdot\}$

Algorithm 25 Derivation mapping for $\mathfrak{p}^{(s)} \xrightarrow{\mathsf{fold}_{\mathcal{S} \to r}^{d,j}} \mathfrak{p}^{(t)}$.

- 1. \triangleright The comments at the top of Algorithm 24 apply here as well.
- 2. **def** $\phi\left(\delta^{(s)}\right)$:
- 3. **case** constant or builtin: **return** $\delta^{(s)}$
- 4. **case** $\bigwedge^{y_1} \circ \cdots \circ \bigwedge^{y_K}$: **return** $\phi \left(\bigwedge^{y_1} \right) \circ \cdots \circ \phi \left(\bigwedge^{y_K} \right)$
- 5. **case** $\begin{pmatrix} x \\ | \\ \downarrow \\ \triangle \end{pmatrix}$ **and** its top rule $\in \mathfrak{p}^{(t)}$: **return** $\phi \in \mathfrak{p}^{(t)}$
- 6. **case** $\begin{pmatrix} x \\ \alpha & \mu & \beta \\ A & A \end{pmatrix}$ **and** its top rule $\notin \mathfrak{p}^{(t)}$:
- 7. return



ightharpoonup The rule passed to u^{-1} is $\mathbf{S}^{(\ell)}$ the top rule of $\boldsymbol{\delta}^{(s)}$ (abusing notation), and the alignment (σ) used in u^{-1} is from the fold transformation's generalized reversibility test.

 \triangleright The rule added is an instantiation of $\mathbf{D}^{(\ell)}$

 \triangleright Solve $\mathbf{T}^{(\ell)} = \theta((x' \oplus \mathbf{P} \alpha' \circ \mu'' \circ \beta'))$

- 8. \triangleright Inverse of the pivoting functions u. See Lemma 3 for formal guarantees.
- 9. **def** $u^{-1}(\delta_2^{(s)}, \mathbf{T}^{(\ell)})$:
- \triangleright Note: With apologies, the argument name $\mathbf{T}^{(\ell)}$ was chosen to match those of u. (It is not a mistake that we pass in a rule from $\mathfrak{p}^{(s)}$ for a variable named $\mathbf{T}^{(\ell)}$ when we call it from the derivation mapping for fold.)
- 10. $\boldsymbol{\delta}_2^{(d)} \leftarrow \boldsymbol{\phi}_1\left(\boldsymbol{\delta}_2^{(s)}\right)$
- 11. $\mathbf{D}^{(\ell)} \leftarrow \overset{s \to a}{\sigma} (\mathbf{T}^{(\ell)})$
- 12. $oldsymbol{y}_j \leftarrow \operatorname{reconstruct}ig(\mathbf{T}^{(\ell)}, \mathbf{D}^{(\ell)}ig)$
- 14. $\boldsymbol{\delta}_1^{(s)} \leftarrow \boldsymbol{\phi}_1 \left(\boldsymbol{\delta}_1^{(d)} \right)$
- 15. **return** $\delta_1^{(s)}$
- 16. \triangleright Reconstruct the intermediate item that would have been erased by $\mathfrak{p}^{(s)} \xrightarrow{\mathsf{unfold} \frac{d,j}{r \to \mathcal{S}}} \mathfrak{p}^{(t)}$.
- 17. **def** reconstruct($\mathbf{T}^{(\ell)}, \mathbf{D}^{(\ell)}$):
- 18. $(x' \oplus = \alpha' \circ y' \circ \beta') \leftarrow r$
- 19. $(y'' \oplus = \mu'') \leftarrow \mathbf{D}^{(\ell)}$
- 20. $\theta \leftarrow \operatorname{cover}((x' \oplus = \alpha' \circ \mu'' \circ \beta'), \mathbf{T}^{(\ell)})$
- 21. return $\theta(y')$

7.4.2.4 Practical Sufficient Conditions for Correctness

This section will identify a practical way to test for the impractical sufficient conditions we identified in the theorems of §7.4.2.3. What will result is a system that can guarantee termination of $\phi(\delta)$ for all derivations δ (where ϕ is either Algorithm 24 or Algorithm 25). We will use a simple abstract interpretation technique to ensure that each recursive call made in ϕ descends in a certain well-founded order—ensuring that the recursive calls eventually hit a base case.

Our strategy for avoiding unsafe fold and unfold transformations is adapted from Roychoudhury et al. (1999). At a high level, the strategy keeps track of the provenance of each rule (i.e., how each rule was derived using previous transformations). The rule provenance data will give rise to additional constraints on folding and unfolding actions between rules of different programs to prevent circular definition cycles. These are merely sufficient conditions but general enough to cover most use cases. A condition worth highlighting is that $p^{(s)}$ and $p^{(d)}$ must be related through a series of transformations; they cannot be unrelated programs.

We start with an initial program $\mathfrak{p}^{(0)}$. The initial program's derivations each have a finite size we can use as a well-founded ordering for induction. As we derive new programs from the initial program through transformations, we can use a common *relative* size with respect to $\mathfrak{p}^{(0)}$ to perform induction on derivations from the later programs. The relative size enables the comparison of derivations across programs (such as the derivations resulting from pivoting in

u and u^{-1}). However, we need to ensure that the relative-size measure remains accurate and well-founded in the new program after each fold or unfold.

Modeling Decisions. To make static analysis tractable, we make the following modeling decisions:

- 1 For any generalized fold or unfold, the programs $\mathfrak{p}^{(s)}$, $\mathfrak{p}^{(d)}$, and $\mathfrak{p}^{(t)}$ are related through a series of fold or unfold transformations to an initial program $\mathfrak{p}^{(0)}$.
- The relative size of each derivation will be (approximately) preserved by each transformation. We will define a function Ψ , such that $\Psi\left(\boldsymbol{\delta}^{(s)}\right) = \Psi\left(\boldsymbol{\delta}^{(t)}\right)$ for all $\left\langle\boldsymbol{\delta}^{(s)},\boldsymbol{\delta}^{(t)}\right\rangle$ in the bijection between $\mathfrak{p}^{(s)}\{\cdot\}$ and $\mathfrak{p}^{(t)}\{\cdot\}$.
- 3 The relative size of a derivation $\delta \in \mathfrak{p}\{\cdot\}$ is an additive function Ψ that is parameterized by values associated with each of the program's rules.

We now make 3 more precise; we define a derivation's measure.

Definition 37. Suppose a given program \mathfrak{p} assigns a **rule measure** $\psi(\mathbf{R})$ to each $\mathbf{R} \in \mathfrak{p}.^{171}$ The **derivation measure** $\Psi(\boldsymbol{\delta})$ for $\boldsymbol{\delta} \in \mathfrak{p}\{\cdot\}$ is

$$\Psi(\boldsymbol{\delta}) \stackrel{\text{def}}{=} \begin{cases}
0 & \text{if } \boldsymbol{\delta} \text{ is a constant or builtin} \\
\Psi(\boldsymbol{\delta}_1) + \dots + \Psi(\boldsymbol{\delta}_K) & \text{if } \boldsymbol{\delta} \text{ is a product } \boldsymbol{\delta}_1 \circ \dots \circ \boldsymbol{\delta}_K \\
\psi(\mathbf{R}) + \Psi(\mathsf{body}(\boldsymbol{\delta})) & \text{if } \boldsymbol{\delta} \text{ is a derivation with top rule } \mathbf{R}
\end{cases} (7.1)$$

We will carefully choose the rule measure parameters (3) so that 2 holds.

After illustrating the following example, we will give the precise recipe.

¹⁷¹We assume, for simplicity, that rule measures are integer-valued. However, it is possible to generalize to certain well-ordered commutative groups (Roychoudhury et al., 1999, Def. 1).

Example 39. Consider the following unfold with rule measures provided to the left of each rule:

$$\begin{pmatrix} 458 & \gamma_0 \colon \mathsf{p} \; + = \; \mathsf{q} \; \star \; \mathsf{t} \, . \\ 459 & & & & \\ 460 & \gamma_1 \colon \mathsf{q} \; + = \; \mathsf{r} \, . \\ 461 & \gamma_2 \colon \mathsf{q} \; + = \; \mathsf{s} \, . \end{pmatrix} \qquad \underbrace{ \begin{array}{c} 462 & \gamma_0 \; + \; \gamma_1 \colon \mathsf{p} \; + = \; \mathsf{r} \; \star \; \mathsf{t} \, . \\ 463 & \gamma_0 \; + \; \gamma_2 \colon \mathsf{p} \; + = \; \mathsf{s} \; \star \; \mathsf{t} \, . \\ 464 & \gamma_1 \colon \mathsf{q} \; + = \; \mathsf{r} \, . \\ 465 & \gamma_2 \colon \mathsf{q} \; + = \; \mathsf{s} \, . \\ 465 & \gamma_2 \colon \mathsf{q} \; + = \; \mathsf{s} \, . \\ 465 & \gamma_2 \colon \mathsf{q} \; + = \; \mathsf{s} \, . \\ 465 & \gamma_2 \colon \mathsf{q} \; + = \; \mathsf{s} \, . \end{pmatrix}$$

We have associated the variables γ_0 , γ_1 , and γ_2 with the 3 rules of the program on the left. The rule measures on the right were chosen to preserve the derivation measure (2). Consider the derivations of 2 derivations for p in each program:

$$\left\{ \begin{array}{ccc} p & p \\ / \backslash & / \backslash \\ q & t & q & t \\ | & | \\ r & s \end{array} \right\} \quad and \quad \left\{ \begin{array}{ccc} p & p \\ / \backslash & / \backslash \\ r & t & s & t \end{array} \right\}$$

We verify that the derivation measures are preserved:

Intuitively, the scheme works because the rules arising from the unfold count for two rules each time they appear in a derivation on the right.

We now consider fold, which works in the opposite direction. Here the measure for each rule appears to its left.

The measure of $^{\#}473$ is interesting: we are uncertain about whether it should be

 $\gamma_3 - \gamma_1$ or $\gamma_2 - \gamma_0$. So, we represent the measure of the new rule using an uncertainty set: $\{\gamma_3 - \gamma_1, \gamma_2 - \gamma_0\}$. We do this because we have lost the information about which of the two rules was folded. Now, we can conservatively estimate the size as an uncertainty set for each pair of derivations in the bijection between

$$\left\{\begin{array}{ccc} p & p \\ / \backslash & / \backslash \\ r & t & s & t \end{array}\right\} \quad and \quad \left\{\begin{array}{ccc} p & p \\ / \backslash & / \backslash \\ q & t & q & t \\ | & | \\ r & s \end{array}\right\}$$

We verify that the derivation measures are preserved (i.e., they are contained within the uncertainty sets):

$$\Psi\left(\begin{array}{c}\mathsf{p}\\ /\backslash\\\mathsf{r}\end{aligned}\right) = \{\gamma_2\} \subseteq \{\gamma_3 - \gamma_1 + \gamma_0, \gamma_2\} = \{\gamma_3 - \gamma_1, \gamma_2 - \gamma_0\} + \{\gamma_0\} = \Psi\left(\begin{array}{c}\mathsf{p}\\ /\backslash\\\mathsf{q}\end{aligned}\right)$$

$$\Psi\left(\begin{array}{c}\mathsf{p}\\ /\backslash\\\mathsf{s}&\mathsf{t}\end{array}\right) = \{\gamma_3\} \subseteq \{\gamma_3,\gamma_2-\gamma_0+\gamma_1\} = \{\gamma_3-\gamma_1,\gamma_2-\gamma_0\} + \{\gamma_1\} = \Psi\left(\begin{array}{c}\mathsf{p}\\ /\backslash\\\mathsf{q}&\mathsf{t}\\ |\\\mathsf{s}\end{array}\right)$$

We make 1 more precise by defining a transformation sequence:

Definition 38. A transformation sequence $\mathfrak{p}^{(0)}, \ldots, \mathfrak{p}^{(T)}$ is defined recursively. Base case (T=0): We allow $\mathfrak{p}^{(0)}$ to be any program. We call $\mathfrak{p}^{(0)}$ the **reference program**. Recursive case $(T\geq 1)$: Suppose $\mathfrak{p}^{(0)}, \ldots, \mathfrak{p}^{(T-1)}$ is a transformation sequence, then $\mathfrak{p}^{(0)}, \ldots, \mathfrak{p}^{(T-1)}, \mathfrak{p}^{(T)}$ is a transformation sequence as long as $\mathfrak{p}^{(T)} = \mathsf{fold}_{S \to r}^{d,j}(\mathfrak{p}^{(s)})$ or $\mathfrak{p}^{(T)} = \mathsf{unfold}_{r \to S}^{d,j}(\mathfrak{p}^{(s)})$ where $\mathfrak{p}^{(s)}, \mathfrak{p}^{(d)}$ are earlier programs in the sequence, the new transformation is well-formed $(\langle S, j, r \rangle)$ are valid arguments), and the new transformation is safe.

Definition 39 (Rule measure updates). For each program in a transformation sequence $\mathfrak{p}^{(t)}$, we define a rule measure $\psi_t(\mathbf{R})$ for each of its rules $\mathbf{R} \in \mathfrak{p}^{(t)}$. The value of the rule measure will generally be an uncertainty set. The rule measure is defined inductively.

Base case: For each rule $\mathbf{R} \in \mathfrak{p}^{(0)}$, we take its measure $\psi_0(\mathbf{R}) > 0$ as an input parameter.¹⁷²

Recursive case: Suppose the last program in our sequence is $\mathfrak{p}^{(t)}$. Then it must have been derived from one of the following cases:

• Case $\mathfrak{p}^{(t)} = \mathsf{unfold}_{r o \mathcal{S}}^{d,j}(\mathfrak{p}^{(s)})$: Here, we have an alignment between each new rule $\mathbf{T}^{(\ell)} \in \mathfrak{p}^{(t)}$ and some $\mathbf{D}^{(\ell)} \in \mathfrak{p}^{(d)}$. That alignment is $\left\{\left\langle \mathbf{D}^{(\ell)}, \mathbf{T}^{(\ell)} \right\rangle\right\}_{\ell=1}^L = \sigma$. The rule measures for each new rule $\mathbf{T}^{(\ell)}$ is then

$$\psi_t(\mathbf{T}^{(\ell)}) = \psi_s(r) + \psi_d(\mathbf{D}^{(\ell)}) \quad \text{for } \ell \in [1:L]$$

• Case $\mathfrak{p}^{(t)} = \operatorname{fold}_{\mathcal{S} \to r}^{d,j}(\mathfrak{p}^{(s)})$: Let r be the new rule that results from folding $\mathcal{S} = \left\{\mathbf{S}^{(\ell)}\right\}_{\ell=1}^L$ by the definitions $\left\{\mathbf{D}^{(\ell)}\right\}_{\ell=1}^L$. Then, r has the following rule measure 173,174

$$\psi_t(r) = \bigcup_{\ell=1}^L (\psi_s(\mathbf{S}^{(\ell)}) - \psi_d(\mathbf{D}^{(\ell)}))$$

¹⁷²The choice of $\psi_0(\mathbf{R}) = 1$ for all \mathbf{R} is fine; however, a clever choice may make certain transformations available and others not. We will use a symbolic system to infer a value for the initial measure that designates a transformation as safe if such a value exists.

¹⁷³The alignment ℓ that we are using here is given by σ of the alignment of the unfold transformation that is used in the generalized reversibility test.

¹⁷⁴Notice that when we fold *multiple* rules |S| > 1, we might introduce approximation since multiple uncertainty sets are merged into $\psi_t(r)$.

In both recursive cases, any unchanged rule $\mathbf{R} \in (\mathfrak{p}^{(t)} \cap \mathfrak{p}^{(s)})$ copies its ψ from $\mathfrak{p}^{(s)}$.

The rule measures for rules in $\mathfrak{p}^{(s)}$ and $\mathfrak{p}^{(d)}$ are defined inductively.

Using relative size to ensure termination. We will devise a simple test that ensures each recursive call of $\underset{s \to t}{\phi}$ descends (i.e., decreases the relative size). We will consider the fold and unfold cases separately.

Consider the case of fold. In this case, we will have an alignment of the form

$$\left\{\left\langle \mathbf{S}^{(\ell)}, \mathbf{D}^{(\ell)} \right\rangle \right\}_{\ell=1}^L = \left\{\left\langle \left(\pmb{x}^{(\ell)} \oplus = \pmb{\alpha}^{(\ell)} \circ \pmb{\mu}^{(\ell)} \circ \pmb{\beta}^{(\ell)} \right), \left(\pmb{y}'^{(\ell)} \oplus = \pmb{\mu}'^{(\ell)} \right) \right\rangle \right\}_{\ell=1}^L$$

We now inspect Algorithm 25's recursive calls. Since the other cases are straightforward, we will focus on line 25.7. Below, we drop the ℓ annotations to avoid notational clutter. In this case, we make the recursive call

$$\phi_{s \to t} \left(\delta^{(s)} \right) = \phi_{s \to t} \left(\begin{array}{c} x \\ \downarrow \\ \triangle \\ \triangle \\ \triangle \\ \end{array} \right) = \phi_{s \to t} \left(\begin{array}{c} x \\ \downarrow \\ \triangle \\ A \\ \end{array} \right) \phi_{s \to t} \left(\begin{array}{c} x \\ \downarrow \\ \triangle \\ \end{array} \right) \phi_{s \to t} \left(\begin{array}{c} x \\ \downarrow \\ \triangle \\ \end{array} \right) \phi_{s \to t} \left(\begin{array}{c} \beta \\ \triangle \\ \triangle \\ \end{array} \right)$$

The tricky thing about these recursive calls is that we need to make sure that the recursive call to $\phi \left(u^{-1} \begin{pmatrix} \mu \\ \triangle \end{pmatrix} \right)$ will not end up calling $\phi \left(\delta^{(s)} \right)$ again. The way that we will operationalize this quest is to use the derivation measure that we developed earlier in this section. Assuming that each derivation has equal measure (2) means that we can focus our attention on proving that $\Psi \left(\delta^{(s)} \right) > \Psi \left(u^{-1} \begin{pmatrix} \mu \\ \triangle \end{pmatrix} \right)$ or, equivalently, $\Psi \left(\delta^{(s)} \right) - \Psi \left(u^{-1} \begin{pmatrix} \mu \\ \triangle \end{pmatrix} \right) > 0$. If we can do that, there is no way that the $\phi _{s \to t}$ program could recurse forever: each step decreases by some nonzero positive integer amount, and the derivation measure is ≥ 0 . It is straightforward to verify that the sibling calls, $\phi _{s \to t} \begin{pmatrix} \alpha \\ \Delta \end{pmatrix}$

and $\phi \left(\bigwedge^{\beta} \right)$, are descending because they are subtrees of $\delta^{(s)}$.

We will now derive the test. First, we expand the derivation measure of $\delta^{(s)}$ by Definition 37:

$$\Psi\left(\boldsymbol{\delta}^{(s)}\right) = \Psi\left(\begin{array}{c} \boldsymbol{x} \\ \boldsymbol{\alpha} & \boldsymbol{\mu} & \boldsymbol{\beta} \\ \boldsymbol{\triangle} & \boldsymbol{\triangle} & \boldsymbol{\triangle} \end{array}\right) = \psi_s(\mathbf{S}) + \Psi\left(\begin{array}{c} \boldsymbol{\alpha} \\ \boldsymbol{\triangle} \end{array}\right) + \Psi\left(\begin{array}{c} \boldsymbol{\mu} \\ \boldsymbol{\triangle} \end{array}\right) + \Psi\left(\begin{array}{c} \boldsymbol{\beta} \\ \boldsymbol{\triangle} \end{array}\right)$$

Recall from Algorithm 25 that

$$u^{-1} \begin{pmatrix} \mu \\ \triangle \end{pmatrix} = \phi \begin{pmatrix} y \\ | \\ \phi \\ s \to d \end{pmatrix}$$

where the rule $\mathbf{D}=(y\oplus \mathbf{P}\mu)$ is an instantiation of a rule in $\mathfrak{p}^{(d)}$. Since this operation is measure-preserving (by assumption), we know that the measure of the new derivation is $\Psi\left(u^{-1}\begin{pmatrix}\mu\\ \triangle\end{pmatrix}\right)=\psi_d(\mathbf{D})+\Psi\begin{pmatrix}\mu\\ \triangle\end{pmatrix}$.

Now, we can put these pieces together:

$$\begin{split} \Psi\left(\boldsymbol{\delta}^{(s)}\right) &= \psi_{s}(\mathbf{S}) + \Psi\left(\stackrel{\alpha}{\triangle}\right) + \Psi\left(\stackrel{\mu}{\triangle}\right) + \Psi\left(\stackrel{\beta}{\triangle}\right) \\ &= \psi_{s}(\mathbf{S}) + \Psi\left(\stackrel{\alpha}{\triangle}\right) + \left(\Psi\left(u^{-1}\left(\stackrel{\mu}{\triangle}\right)\right) - \psi_{d}(\mathbf{D})\right) + \Psi\left(\stackrel{\beta}{\triangle}\right) \\ \Psi\left(\boldsymbol{\delta}^{(s)}\right) - \Psi\left(u^{-1}\left(\stackrel{\mu}{\triangle}\right)\right) &= \psi_{s}(\mathbf{S}) + \Psi\left(\stackrel{\alpha}{\triangle}\right) - \psi_{d}(\mathbf{D}) + \Psi\left(\stackrel{\beta}{\triangle}\right) \end{split}$$

Thus, we have descent if the following holds for all possible siblings $\stackrel{\alpha}{\bigtriangleup}$ and $\stackrel{\beta}{\bigtriangleup}$

and all possible measures in the uncertainty set:^{175,176}

$$(\psi_s(\mathbf{S}) - \psi_d(\mathbf{D})) + \Psi\begin{pmatrix} \alpha \\ \triangle \end{pmatrix} + \Psi\begin{pmatrix} \beta \\ \triangle \end{pmatrix} > 0$$

In practice, we replace $\Psi\left(\stackrel{\alpha}{\triangle} \right)$ and $\Psi\left(\stackrel{\beta}{\triangle} \right)$ with estimates $\widehat{\Psi}(\alpha)$ and $\widehat{\Psi}(\beta)$ that only depend on the top of the derivation. These estimates must satisfy

$$\widehat{\Psi}(\mathbf{x}) \leq \min_{\mathbf{\delta} \in \mathfrak{p}^{(s)}\{\mathbf{x}\}} \Psi(\mathbf{\delta}) \quad ext{for all } \mathbf{x} \in \mathbb{H}$$

Notice that measure preservation means that we can use any program's derivation set for the base of the minimization. Thus, $\widehat{\Psi}$ can be estimated from the initial program $\mathfrak{p}^{(0)}$, which has the additional benefit that its derivation measure has the highest fidelity.¹⁷⁸ We describe possible estimators in Definition 40.

For the case of the unfold transformation, we arrive through a similar argument at the following condition:¹⁷⁹

$$\psi_s(\mathbf{S}) + \psi_d(\mathbf{D}) + \Psi\begin{pmatrix} \alpha \\ \triangle \end{pmatrix} + \Psi\begin{pmatrix} \beta \\ \triangle \end{pmatrix} > 0$$

Notice that $\psi_s(\mathbf{S}) + \psi_d(\mathbf{D})$ need not be positive.

Lastly, we must ensure the entire provenance chain of $\mathfrak{p}^{(t)}$ is safe. This means that as we crawl up the sequence of prior fold and unfold transformations, we need to check that each is safe and they derive from a common $\mathfrak{p}^{(0)}$.

¹⁷⁵When x is an uncertainty set, the stipulation x > 0 is shorthand for $\forall x_i \in x : x_i > 0$.

¹⁷⁶This condition is equivalent to that of Roychoudhury et al. (1999). However, we have arrived at it through different arguments.

¹⁷⁷Notice that our choice to use the relative size of a derivation rather than relative height pays off here because the siblings in the derivation can provide some slack for the condition.

¹⁷⁸Its uncertainty sets ought to be singletons.

¹⁷⁹This condition was missing in Roychoudhury et al. (1999) since they had not considered generalized unfold.

Definition 40. The following are possible estimators for the **minimum derivation measure** of an item.

- The simplest estimator is $\widehat{\Psi}(\cdot) = 0$.
- A simple estimator is

$$\widehat{\Psi}(oldsymbol{x}) \stackrel{ ext{def}}{=} egin{cases} 0 & ext{\it if x is a constant, builtin, or input} \ \widehat{\Psi}(oldsymbol{x}_1) + \dots + \widehat{\Psi}(oldsymbol{x}_K) & ext{\it if x is a product $x_1 \circ \dots \circ x_K$} \ 1 & ext{\it otherwise} \end{cases}$$

• A more precise $\widehat{\Psi}$ can be expressed through the following system of linear inequalities that are a sufficient characterization of the set of feasible estimators:¹⁸⁰

$$\begin{cases} &\widehat{\Psi}(\cdot)\geq 0\\ &\textit{for } \textit{x in constants, builtins, or inputs}: \ \widehat{\Psi}(\textit{x})\leq 0\\ &\textit{for } \mathbf{R}\in \mathfrak{p}^{(0)}:\\ &(\textit{x} \oplus = \textit{y}_1 \circ \cdots \circ \textit{y}_K) \leftarrow \mathbf{R}\\ &\widehat{\Psi}(\textit{x}) \leq \psi_0(\mathbf{R}) + \widehat{\Psi}(\textit{y}_1) + \cdots + \widehat{\Psi}(\textit{y}_K) \end{cases}$$

For $\mathbf{R} \in \mathfrak{p}^{(0)}$, $\psi_0(\mathbf{R})$ should be a singleton set; thus, we can treat it like a scalar in these inequalities.

We summarize the complete system in Algorithm 26.

Theorem 1 (Correctness of Algorithm 26). Let $\mathfrak{p}^{(0)}, \mathfrak{p}^{(1)}, \ldots, \mathfrak{p}^{(t)}$ be a (measure-safe) transformation sequence. Then, we can assess the safety of a candidate fold or unfold transformation, $\mathfrak{p}^{(t)} \to \mathfrak{p}^{(?)}$, using measure_safe($\mathfrak{p}^{(?)}$). If it is true, then the transformation is safe.

¹⁸⁰ For efficiency, we might coarsen the program, e.g., by a maximum depth, as we did in §6.1.

Proof Sketch. The measure safety test (Algorithm 26) ensures that the recursive procedures for derivation mappings (Algorithm 24 and Algorithm 25) always terminate. The method measure_safe($\mathfrak{p}^{(?)}$) will only return true when the recursive calls provably descend safely under the candidate fold or unfold. The argument for why this is the case was detailed earlier in this section. Therefore, the candidate transformation must be safe by Theorem 2 and Theorem 3.

This section has provided an effective safety test for the Summary. weighted setting based on prior work by Roychoudhury et al. (1999) for the unweighted setting. We have found that this system provides excellent recall when measure safe leverages a constraint solver. 181 Not requiring an initial rule measure specification makes the overall system much more practical.

¹⁸¹We have found the complicated stratified counters in Roychoudhury et al. (1999) to be unnecessary thanks to the constraint solver.

Algorithm 26 Fold–unfold transformation sequence safety test. We assume an initial reference program $\mathfrak{p}^{(0)}$ with rule measure $\psi_0(\cdot)$ as input.

```
\widehat{\Psi}(\cdot) \leftarrow \{ \text{one of the estimators in Definition 40} \}
 2. def \psi_t(\mathbf{R}):
           case \mathfrak{p}^{(t)}=\mathfrak{p}^{(0)}: return \psi_0(\mathbf{R})
          case \mathfrak{p}^{(t)} = \text{fold}_{S \to r}^{d,j}(\mathfrak{p}^{(s)}) with alignment \sigma:
                    return \bigcup_{\langle \mathbf{S}, \mathbf{D} \rangle \in \sigma} (\psi_s(\mathbf{S}) - \psi_d(\mathbf{D}))
               else
                    return \psi_{s}(\mathbf{R})
                                                                                                                                  ▷ Copy the unchanged rule
          case \mathfrak{p}^{(t)} = \mathsf{unfold}_{r \to \mathcal{S}}^{d,j}(\mathfrak{p}^{(s)}) with alignment \sigma:
10
                    return \psi_s(r) + \psi_d \left( \sigma(\mathbf{R}) \right)
               else
12
                    return \psi_{\mathfrak{s}}(\mathbf{R})
                                                                                                                                ▷ Copy the unchanged rules
13.
    \mathbf{def}\ \varphi(t):
          case \mathfrak{p}^{(t)} = \mathfrak{p}^{(0)}: return true
          case \mathfrak{p}^{(t)} = \text{fold}_{S \to r}^{d,j}(\mathfrak{p}^{(s)}) with alignment \sigma:
16
               	riangleright In the code below, S is of the form <math>(x \oplus \exists \alpha \circ \mu \circ \beta) where \alpha, \mu, and \beta are deter-
                   mined by how S was folded.
              \mathbf{return}\ \varphi(s) \land \varphi(d) \land \bigwedge_{\langle \mathbf{S}, \mathbf{D} \rangle \in \sigma} \left( \psi_s(\mathbf{S}) - \psi_d(\mathbf{D}) + \widehat{\Psi}(\boldsymbol{\alpha}) + \widehat{\Psi}(\boldsymbol{\beta}) > 0 \right)
18
          \mathbf{case}\; \mathfrak{p}^{(t)} = \mathsf{unfold}_{\; r \to \mathcal{S}}^{\; d,j}(\mathfrak{p}^{(s)}) \; \mathbf{with \; alignment} \; \sigma\text{:}
               (\mathbf{x} \oplus \mathbf{=} \boldsymbol{\alpha} \circ \mathbf{y}_{j} \circ \boldsymbol{\beta}) \leftarrow r \qquad \triangleright \text{Here } \alpha, \ y_{j}, \ \text{and } \beta \text{ are determined by how } r \text{ was unfolded.} 
 \mathbf{return} \ \varphi(s) \land \varphi(d) \land \bigwedge_{\langle \mathbf{D}, \mathbf{T} \rangle \in \sigma} \left( \psi_{s}(\mathbf{S}) + \psi_{d}(\mathbf{D}) + \widehat{\Psi}(\boldsymbol{\alpha}) + \widehat{\Psi}(\boldsymbol{\beta}) > 0 \right) 
20
21
22. def measure_safe(\mathfrak{p}^{(t)}):
          \triangleright For best results, use a constraint solver to find values for \{\gamma_{\mathbf{R}}\}_{\mathbf{R}\in\mathfrak{p}^{(0)}} that satisfy
              the proposition \varphi(t). Recall that each \gamma_{\mathbf{R}} must be > 0 and an integer. Our imple-
              mentation uses the popular Z3 solver (Moura and Bjørner, 2008), but we note
              that any integer linear solver would suffice.
          \triangleright (For faster results, we can check the feasibility of \varphi(t) using \gamma_{\mathbf{R}}=1 for all \mathbf{R}\in
24
          (It is always ok to return false, e.g., if we want to set a timeout on the SMT
25
          return true if \varphi(t) is satisfiable else false
26
```

7.5 Hoisting Transformations

This section gives two related transformations: the speculation transformation (Eisner and Blatz, 2007) and the left-corner transformations (Rosenkrantz and Lewis, 1970; Johnson, 1998; Johnson and Roark, 2000). Our left-corner transformation generalizes the well-known left-corner transformation for context-free grammars to Dyna programs. We find an interesting connection between the left-corner and speculation transformations: both transformations "hoist" a subgoal up the tree, but the left-corner transformation additionally *transposes* the structure of each derivation—replacing left recursion with right recursion.

7.5.1 Speculation Transformation

The speculation transformation was first introduced in Eisner and Blatz (2007). Our presentation will differ from theirs because we will describe it in terms of how it transforms derivations.

Example 40. Consider again the CKY program (Example 6). 183

```
474 p(X,I,K) += p(Y,I,J) * p(Z,J,K) * r(X,Y,Z).

475 p(X,I,K) += p(Y,I,K) * r(X,Y).

476 p(X,I,K) += w(Y,I,K) * r(X,Y).
```

This program has an inefficiency: it repeats the same unary-chain closure for all spans $\langle I, K \rangle$. Our use of speculation will refactor the program so that unary chain closures may be precomputed (independently of the sentence) and reused at

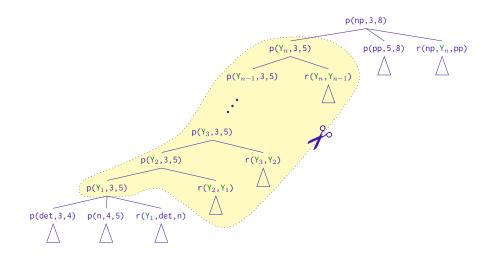
¹⁸²The left-corner transformation eliminates left recursion in context-free grammars, which makes them amenable to specialized top-down parsing algorithms.

¹⁸³The following example is adapted from Eisner and Blatz (2007).

each span. This trick—attributed to Stolcke (1995)—is standard in probabilistic parsing. The final program is as follows:

```
 u(X,X') += u(Y,X') * r(X,Y). 
 u(X,X') += r(X,X'). 
 p(X,I,K) += p'(X',I,K) * u(X,X'). 
 480  p(X,I,K) += p'(X,I,K). 
 p'(X,I,K) += p(Y,I,J) * p(Z,J,K) * r(X,Y,Z). 
 481  p'(X,I,K) += w(Y,I,K) * r(X,Y).
```

But how were we able to arrive at it? Below is a derivation of the original program that contains a long chain of #475 all with 1=3, K=5 (highlighted):

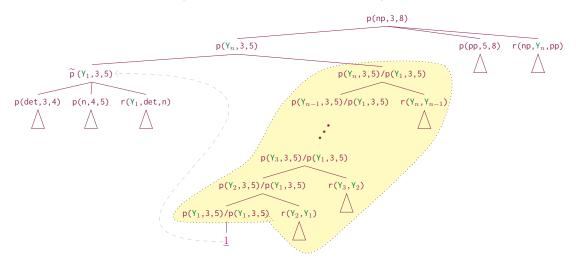


Our transformed program will use a slashed item¹⁸⁴ $p(Y_n,3,5)/p(Y_1,3,5)$ that computes the total weight of the intermediate chains—specifically, those that start at $p(Y_1,3,5)$, end at $p(Y_n,3,5)$, and go through a subset of rules (in this case $\{\#475\}$). ¹⁸⁵ These different chains will correspond to different derivations of the slashed item. We can think of the new item as abstracting away the chain from the specific starting point $p(Y_1,3,5)$. We can commit to a specific starting point of the form $p(Y_1,3,5)$ later (i.e., higher up in the derivation tree) using a recovery rule, such as

¹⁸⁴There is nothing special about the use of / here—it is not a special operator; it is just a name for the new item that is meant to *suggest* that $p(Y_n, 3, 5) = p(Y_1, 3, 5) * p(Y_n, 3, 5)/p(Y_1, 3, 5)$ as it would if / meant division under real arithmetic.

¹⁸⁵The specific subset of rules is a parameter of the transformation.

 $p(Y_n,3,5) += \widetilde{p}(Y_1,3,5) * p(Y_n,3,5)/p(Y_1,3,5)$. For technical reasons, we rename the bottom element of the chain apart from its original by marking it with $\widetilde{\cdot}$ (in this case, $\widetilde{p}(Y_1,3,5)$). Our transformed derivation looks as follows.



In the diagram, we can see that the chain has been stashed into a derivation of the slashed item. That subderivation turns out not to depend on the specific values for the span $\langle 3,5 \rangle$ —meaning that it can be reused by different spans and leads to a runtime improvement after we apply abbreviation (§6.2.3). The subderivations represented as generic triangles \wedge will be transformed recursively.

In the speculation notation, our transformed program is

```
483 p(X,I,K) / p(X',I',K') += p(Y,I,K) / p(X',I',K') * r(X,Y).

484 p(X',I',K') / p(X',I',K').

485 p(X,I,K) += \widetilde{p}(X',I,K) * p(X,I,K) / p(X',I',K').

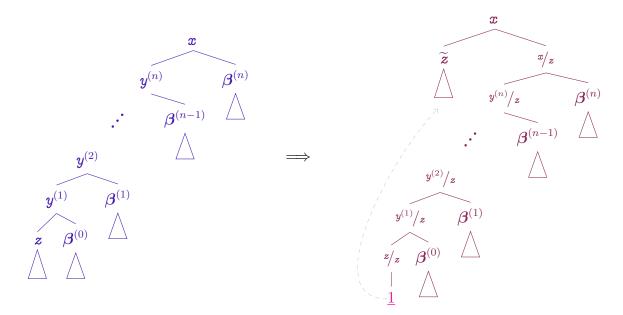
486 \widetilde{p}(X,I,K) += p(Y,I,J) * p(Z,J,K) * r(X,Y,Z).

487 \widetilde{p}(X,I,K) += w(Y,I,K) * r(X,Y).
```

Program specialization (§6.2) recognizes that I=I' and K=K'. Additionally, it recognizes that I, I', K, and K' are always free in the slashed items so they can be dropped. So after specialization, we have the final program from earlier.

The speculation transformation is parameterized by $\sigma \subseteq \mathfrak{p}$ and a nonground term z. Let $\delta \in \mathfrak{p}\{\cdot\}$ be an arbitrary derivation tree of a program \mathfrak{p} . We may

represent δ schematically as a derivation of the form on the left. The speculation transform will transform δ into a derivation of the form on the right.



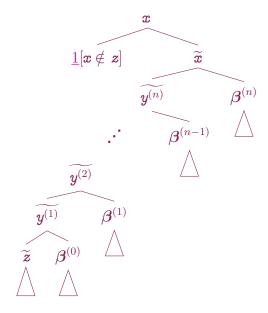
The schematic above exposes a **spine** of the derivation: the (maximal length) path $x \rightsquigarrow z$ along the left edge of the derivation that uses only some designated subset of rules σ and ends at an item matching the nonground term z.¹⁸⁶ As we can see with the derivation on the right, the speculation transformation has $dragged\ up$ (hoisted) the subderivation rooted at z leaving behind a **trace** $x/z \rightsquigarrow z/z$ along the spine $x \rightsquigarrow z$ with a 1 in the original position of the hoisted subderivation. We call each 1/z item a **slashed item**. The new derivation still has root x, but now z attaches at the top with a derivation of the x item without z at the base of its spine. Each subtree rooted at z0 remains a derivation of the same z0 items, but the internals of each subtree may change because the transformation is applied recursively to it.

Recall from Example 40 that the marking \tilde{z} on the z item means that the item

 $[\]overline{^{186}}$ Thus, if z appears multiple times, we take the bottom-most occurrence.

is the bottom of a chain. We will use this marking for an additional purpose. In \mathfrak{p} , it might also be possible to derive x in other ways using the rules in σ such that the base of the spine does not match z. We want to keep these derivations intact, so we create a route marked with $\tilde{\cdot}$ that allows the derivations built with rules in σ that do not have z along their spine to continue contributing values to x.

These marked items will contribute to the value of x using a rule $(x \oplus = \widetilde{x})$, but notice that we need to ensure that $x \notin z$ because the \widetilde{x} is supposed to sum up the derivations of x that either do not have z on their spine—either because (a) they do not go through σ , or (b) they do go through σ , but they do not end in a z item. If x is a z item, it reaches z with a zero-length spine. To drop these unwarranted contributions, Eisner and Blatz (2007) add a constraint that $x \notin z$.



The last case we need to consider is when multiple z items might appear along the spine. If we do not distinguish this case, we could overcount, as there are many places where the \tilde{z} item will try to attach itself. We use the solution proposed by Eisner and Blatz (2007): hoist the bottom-most z item as high as

it can go along the spine (i.e., breaking the ambiguity in where to attach it by attaching at the highest point along the spine).

Definition 41. Given a program \mathfrak{p} , the **speculation transformation** $\mathfrak{p}'=$ speculation(\mathfrak{p}, z, σ) is parameterized by a nonground term z representing a set of items to slash out, and σ is a set of rule indices to slash from. For simplicity, the transformation given below requires that (1) \mathfrak{p} does not have inputs, $\mathfrak{1}^{187}$ (2) the variables in z cannot appear in \mathfrak{p} , (3) for each $i \in \sigma$, $(x \oplus y_1 \circ \cdots \circ y_K) = \mathfrak{p}_i$:

(a) K > 0, (b) y_1 is not a constant or builtin, $\mathfrak{1}^{188}$ (c) $y_1 \subseteq z$ or $y_1 \cap z = \emptyset$. $\mathfrak{1}^{189}$

Algorithm 27 speculation: Apply the speculation transformation (Definition 41).

```
1. def speculation(\mathfrak{p}, z, \sigma):
            \mathfrak{p}' \leftarrow 7
           for i = 1 ... |p|:
                  if i \notin \sigma:
                        (x \oplus = \mu) \leftarrow \mathfrak{p}_i
                        \mathfrak{p}'[\widetilde{x}] \oplus = \mu
                        (\boldsymbol{x} \oplus \boldsymbol{=} \boldsymbol{y} \circ \boldsymbol{\beta}) \leftarrow \mathfrak{p}_i
                        \mathfrak{p}'[x/z] \oplus = y/z \circ \beta
                        if y \cap z = \emptyset:
10
                              \mathfrak{p}'[\widetilde{x}] \oplus = \widetilde{y} \circ \beta
11.
           \triangleright Base case for slashed items
            \mathfrak{p}'[z/z] \oplus = 1
13.
           \triangleright Rules that recover the original items in \mathfrak{p}
14.
           \mathfrak{p}'[X] \oplus = \widetilde{z} \circ X/z
15.
            \mathfrak{p}'[X] \oplus = \text{not-matches}_{z}(X) \circ \widetilde{X}
           return p'
17.
```

conditions: $(y_1 \cap z)$, $(y_1 \setminus z)$, and $(z \setminus y_1)$.

¹⁸⁷Annoyingly, when the program has inputs, the transformation given will rename the inputs (wrapping them in $\widetilde{}$)—breaking the existing program API. It is straightforward to relax this assumption when the names of the inputs are disjoint from items that are defined by the program: (1) override \widetilde{x} not to mark x when it is an input relation, and (2) constrain line 27.16 so that it will not produce the problematic rule $x \oplus = \operatorname{not-matches}_z(x) \circ x$.

 $^{^{188}}$ To support slashing-out constants and built-in relations, we can wrap them in temporary items, and then eliminate them (by unfolding) once the speculation transformation is complete. 189 It is sometimes useful to partition \mathfrak{p}_i into three more specialized rules using boolean side

Algorithm 28 not-matches: Implementation of the not-matches relation.

 \triangleright This method implements not-matches_z(Y) $\stackrel{\text{def}}{=}$ 1[Y \notin z], a built-in relation that checks if Y is not an element of the set of terms denoted by the nonground term z. Importantly, the variables in z are part of the name of the builtin, so they are not summed over when the builtin is used in a rule. This method should be included in the library of built-in relations. It is implemented as a result stream to be compatible with §3.1.2.2.

def not-matches_z(y):

```
if unify(z, y) = FAIL: yield \underline{1} \triangleright certain, satisfied: z and y cannot match else if cover(z, y) \neq FAIL: return \triangleright certain, fail: z \supseteq y means that y must be in z else yield not-matches<sub>z</sub>(y) \triangleright uncertain, delay: y could still become an instance of z
```

Extensions. We also note that a straightforward adaptation to Definition 41 speculates over the rightmost subgoal instead of the leftmost subgoal. Our implementation provides both leftmost and rightmost options. When the semiring is commutative, we can permute the subgoals (as a pre-transformation) to make any subgoal the first subgoal in its rule. Even without commutativity, some subgoals may commute with one another—in particular, built-in constraints can be permuted because they multiply a 1 or 0, which can be freely moved around thanks to the semiring properties. Lastly, we note that many of

¹⁹⁰The strategy of using a built-in relation is a workaround that avoids a mixed semiring program as was done in Eisner and Blatz (2007)'s presentation.

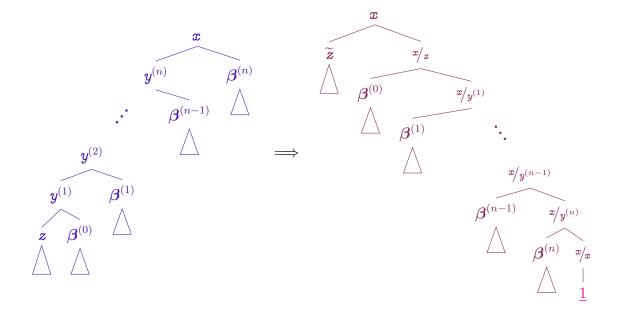
the $\widetilde{\cdot}$ items created by the speculation transformation are unnecessary; thus, it is often worth following speculation with unfolding and useless-rule elimination to remove some unnecessary intermediate $\widetilde{\cdot}$ items. An easy choice is to avoid creating $\widetilde{\cdot}$ rules for any item that is not in $\{z\} \cup \{y \mid (x \oplus \exists y \circ \beta) = \mathfrak{p}_i, i \in \sigma\}$.

Correctness. We believe that the speculation transformation is strongly semantics-preserving. The / and ~ items should be regarded as new definitions (§7.3), as they are not defined in p. We have provided a derivation mapping for the transformation in §B.3. However, we have not provided a complete proof of correctness.

7.5.2 Left-Corner Transformation

We extend the left-corner transformation (Rosenkrantz and Lewis, 1970; Johnson, 1998) to Dyna programs. To have fine-grained control over how the transformation is applied, we build on the selective left-corner transformation (Johnson and Roark, 2000), which permits the transformation to apply to a subset of rules. The left-corner transformation is closely related to speculation: both transformations hoist up the subderivations, but the left-corner transformation will, additionally, switch left recursion to right recursion (i.e., it will transpose the spine of the derivation). Below is a schematic that builds on the earlier discussion of speculation.

¹⁹¹This set of items is useful for specializing the recovery rules on line 27.15 and line 27.16.



In contrast with speculation, the slashed items in the trace of the transformed derivation are named differently; rather than naming them by the item z that is hoisted up, they are named by the destination y. This is how the transformation manages to transpose the derivation tree. This effect can be beneficial when the program defines many items irrelevant to a set of output items. ¹⁹² Our version of the left-corner transformation does not constrain the base of the spine, so there is no need to create the recursive rules of $\tilde{\cdot}$ items. We maintain the $\tilde{\cdot}$ items to distinguish items at the base of the spine from the top.

Definition 42. Given a program \mathfrak{p} , left-corner transformation $\mathfrak{p}' = \mathsf{leftcorner}(\mathfrak{p}, \sigma)$ is parameterized by a set of rule indices σ . For simplicity, the transformation be-

¹⁹²This is analogous to forward-mode vs. backward-mode automatic differentiation (Griewank and Walther, 2008). Specifically, speculation is to left-corner as forward-mode is to backward-mode. The differentiation methods build items of the form $\frac{\partial x}{\partial z}$, much like hoisting methods build items of the form $^x/z$. However, they build different subsets: backward-mode and left-corner target the numerators, and forward-mode and speculation target the denominators. It is well-known in the automatic differentiation setting, that backward-mode is more efficient when the number of numerators needed is less than the number of denominators needed, and forward-mode is more efficient in the opposite regime.

low requires that (1) $\mathfrak p$ does not have inputs, ¹⁹³ (2) for each $i \in \sigma$, $(x \oplus \exists y_1 \circ \cdots \circ y_K) = \mathfrak p_i$:

(a) K > 0, (b) y_1 is not a constant, builtin, or input. ¹⁹⁴

Algorithm 29 leftcorner: Apply the left-corner transformation (Definition 42).

```
1. def leftcorner(\mathfrak{p}, \sigma):
            \mathfrak{p}' \leftarrow 7
            for i = 1 ... |p|:
                  if i \notin \sigma:
                         (\boldsymbol{x} \oplus \boldsymbol{=} \boldsymbol{\mu}) \leftarrow \mathfrak{p}_i
                        \mathfrak{p}'[\widetilde{x}] \oplus = \mu
                  else
                         (\boldsymbol{x} \oplus \boldsymbol{\exists} y \circ \boldsymbol{\beta}) \leftarrow \mathfrak{p}_i
                        \mathfrak{p}'[X/y] \oplus = \beta \circ X/x
            \triangleright Base case of slashed items
10.
            \mathfrak{p}'[X/X] \oplus = 1
           \triangleright Rules that recover the original items in \mathfrak p
            \mathfrak{p}'[X] \oplus = \widetilde{Z} \circ X/Z
           \mathfrak{p}'[X] \oplus = \widetilde{X}
14.
            return p'
15.
```

The comments in the bottom paragraph of Definition 41 apply here.

Extensions. The same extensions for the case of speculation apply here: the right-corner transformation that hoists the rightmost subgoals, and permuting commutative subgoals.

Correctness. We believe that the left-corner transformation is strongly semantics-preserving. It is sufficiently similar to the grammar transformation that we do not foresee any issues. However, we have not formally checked all of the details.

 $[\]overline{^{193}\text{See}}$ footnote 187 for a workaround.

¹⁹⁴See footnote 188 for a workaround.

7.5.3 Examples

Example 41. What happens if we apply the left-corner transformation to Example 40? The transformed program is almost identical: both produce precisely the same items. However, as expected, the left-corner transformation will transpose the recursion: the slashed items are the same but defined using a right-recursive rule instead of a left-recursive rule. Specifically, the new program is the same as #483–487, except that we replace #483 (repeated below)

```
488 p(X,I,K) / p(X',I',K') += p(Y,I,K) / p(X',I',K') * r(X,Y).

with the rule
```

489 p(X',I',K') / p(Y,I,K) += r(X,Y) * p(X',I',K') / p(X,I,K).

Example 42. This example explores the connection between speculation and the variable-elimination strategy we saw in §7.3. Consider the following single-rule program, which may be regarded as the product of three matrices:

```
490 a(I,L) += b(I,J) * c(J,K) * d(K,L).
```

Factoring the last two subgoals into an intermediate item tmp(J,L):

```
a(I,L) += b(I,J) * tmp(J,L).

tmp(J,L) += c(J,K) * d(K,L).
```

The new program is more efficient because a rule's maximum number of variables is reduced. We can mimic this transformation by speculation with $\sigma = \{ \text{\#}490 \}$ and $\mathbf{z} = b(\text{I'},\text{J'})$:

```
493 a(I,L) / b(I',J') += b(I,J) / b(I',J') * c(J,K) * d(K,L).
494 a(I,L) += a(I,L) / b(I',J') * b(I',J').
495 b(I',J') / b(I',J').
```

Next, we eliminate #495 *by unfolding:*

```
496 a(I,L) += a(I,L) / b(I,J) * b(I,J).
497 a(I,L) / b(I,J) += c(J,K) * d(K,L).
```

Notice that the second rule always has I free. Thus, speculation's items a(I,L) / b(I,J) are effectively just a naming scheme for the tmp(J,L) items, a step that we can automate with the abbreviation transformation (§6.2.3).

Example 43. Suppose we have a program such as

```
498  goal += x(I) * f(I).
499  goal += x(I) * g(J) * z.

and we want to factor it into the following form:
500  goal += x(I) * (???).

We can do so by speculation with z = x(I) and σ = {#498, #499}:
501  goal += x(I') * goal/x(I').
502  goal/x(I') += x(I)/x(I') * f(I).
503  goal/x(I') += x(I)/x(I') * g(J).
504  x(I')/x(I').

Next, we unfold the x(I)/x(I') subgoals:
505  goal += x(I') * goal/x(I'). % unchanged
506  goal/x(I') += f(I').
507  goal/x(I') += g(J) * z.
```

Notice that unfolding binds I to I' in #506, and #507 is non-range-restricted (i.e., has a free variable in its head). We see that goal/x(I') is a clever scheme for naming the contributions from both of these rules. This program has the form that we were after.

Example 44. A split bilexical grammar is one where the head word must combine with all its right children before any of its left children. The naive algorithm, given below, runs in $\mathcal{O}(n^5)$. Eisner and Satta (1999a) were able to speed it up to $\mathcal{O}(n^3)$. Eisner and Blatz (2007) described the speedup in terms of speculation. To

simplify the example, we assume an unlabeled grammar.

```
508 r(I,I,K) += word(I,K).  % no right children yet

509 r(P,I,K) += r(P,I,J) * arc(P \rightarrow C) * I(C,J,K). % add right child

510 % no left children yet

511 I(P,I,K) += r(P,I,K). % no left children yet

512 I(P,I,K) += I(C,I,J) * arc(C \leftarrow P) * I(P,J,K). % add left child

513 goal += I(P,\emptyset,N) * len(N).
```

The crux of the speedup is to transform left constituents so they can be built without a specific right constituent. This work can be done independently of the right constituent, and the slashed items we will derive in the transformed program below do not instantiate K—leading to an improved runtime. Below is the output of speculation from the right with z = r(H', I', K') and $\sigma = \{\#511, \#512\}$ after omitting the useless \widetilde{goal} item and useless-rule elimination:

```
515 \widetilde{r}(I,I,K) += word(I,K).

516 \widetilde{r}(P,I,K) += r(P,I,J) * arc(P \rightarrow C) * 1(C,J,K).

517

518 I(P,I,K) / r(H',I',K') += r(P,I,K) / r(H',I',K').

519 I(P,I,K) / r(H',I',K') += 1(C,I,J) * arc(C \leftarrow P) * 1(P,J,K) / r(H',I',K').

520 I(P,I,K) / r(H',I',K') * len(N).

521 I(P,I,K) / r(H',I',K') / r(H',I',K')).

522 I(P,I,K) += I(P,I,K) / r(P,I,K) / r(H',I',K').

523 I(P,I,K) += I(P,I,K) / r(H',I',K') * 1(P,I,K) / r(H',I',K').
```

Next, we eliminate all uses of the original r and 1 items using unfold; and we run program specialization, which exposes variables that are equated:

```
\begin{array}{lll} \widetilde{r}(I,I,K) \; += \; \text{word}(I,K) \, . \\ \\ 530 \;\;\; \widetilde{r}(I,I,K) \; += \; \widetilde{r}(I,I,J) \; * \; \text{arc}(I \; \rightarrow \; \mathbb{C}) \; * \; \widetilde{r}(\mathbb{C},\mathbb{C},K) \; * \; \mathbb{1}(\mathbb{C},J,K) \; / \; r(\mathbb{C},\mathbb{C},K) \, . \\ \\ 531 \;\;\; \mathbb{1}(I',I',K') \; / \; r(I',I',K') \, . \\ \\ 532 \;\;\; \mathbb{1}(I',I,K') \; / \; r(I',I',K') \; += \; \widetilde{r}(\mathbb{C},\mathbb{C},J) \; * \; \mathbb{1}(\mathbb{C},I,J) \; / \; r(\mathbb{C},\mathbb{C},J) \\ \\ & \;\;\; * \;\; \text{arc}(\mathbb{C} \; \leftarrow \; I') \; * \; \mathbb{1}(I',J,K') \; / \; r(I',I',K') \, . \\ \\ 534 \;\;\; \text{goal} \;\; += \; \widetilde{r}(H,H,N) \; * \; \mathbb{1}(H,\emptyset,N) \; / \; r(H,H,N) \; * \; \text{len}(N) \, . \end{array}
```

We see that the first two arguments of $\widetilde{r}(H,H,N)$ are always equal; so we can

 $abbreviate \ \widetilde{r}(H,H,N) \mapsto rtrap(H,N).$

```
7535 rtrap(I,K) += word(I,K).

7536 rtrap(I,K) += rtrap(I,J) * arc(I \rightarrow C) * rtrap(C,K) * l(C,J,K) / r(C,C,K).

7537 l(I',I',K') / r(I',I',K').

7538 l(I',I,K') / r(I',I',K') += rtrap(C,J) * l(C,I,J) / r(C,C,J)

7539 * arc(C \leftarrow I') * l(I',J,K') / r(I',I',K').

7540 goal += rtrap(H,N) * l(H,0,N) / r(H,H,N) * len(N).
```

We chose this name because it corresponds to Eisner and Satta (1999a)'s notion of a "right trapezoid;" we will derive their "left trapezoid" next. First, we observe the redundancy in 1(I',I,K') / r(I',I',K'), which has three variables equated. We can abbreviate $1(I',I,K') / r(I',I',K') \mapsto 1r(I',I,K')$.

```
741 rtrap(I,K) += word(I,K).

542 rtrap(I,K) += rtrap(I,J) * arc(I \rightarrow C) * rtrap(C,K) * lr(C,J,K).

543 lr(I',I',K').

544 lr(I',I,K') += rtrap(C,J) * lr(C,I,J) * arc(C \leftarrow I') * lr(I',J,K').

545 goal += lr(H,0,N) * rtrap(H,N) * len(N).
```

Next, we observe that K' in Ir(I',I,K') is always free, so we drop it by abbreviating

Finally, we apply variable elimination to derive the following new items:

```
tmp(C,I) += ltrap(I,J) * rtrap(C,J).
```

Eagerly folding this new definition into the program gives us our final $O(n^3)$ time algorithm:

```
552  rtrap(I,K) += word(I,K)

553  rtrap(I,K) += tmp(I,C) * arc(I \rightarrow C) * rtrap(C,K).

554  ltrap(I',I').

555  ltrap(I',I) += ltrap(C,I) * arc(C \leftarrow I') * tmp(C,I').

556  tmp(C,I) += ltrap(I,J) * rtrap(C,J).

557  goal += ltrap(H,0) * rtrap(H,N) * len(N).
```

Example 45. Consider the following left-recursive program for computing the total weight of all paths in a graph.

```
startpath(X) += start(X).
startpath(Y) += startpath(X) * edge(X,Y).
goal += startpath(X) * stop(X).
```

If we apply the left-corner transformation to it, we get the following right-recursive program:

```
561 (goal / start(X)) += goal / startpath(X).
562 (goal / startpath(X)) += edge(X,Y) * goal / startpath(Y).
563 (goal / startpath(X)) += stop(X).
564 goal += start(X) * goal / start(X).
```

Next, we unfold goal / start(X) in the bottom rule and abbreviate (goal / startpath(X))

```
stoppath(X) += edge(X,Y) * stoppath(Y).
stoppath(X) += stop(X).
goal += start(X) * stoppath(X).
```

 \mapsto stoppath(X):

This program computes the same value for goal, but it does so using the total weight of paths starting at a given node instead of the total weight of all paths ending at a given node. Both are well-known methods for computing the total weight of all paths. Here, we have presented a mechanical technique for deriving one from the other. (Note that if we apply the same exercise to this last program using the right-corner transformation, we will arrive at the original program.)

Example 46. Recall Example 28. That program (repeated below) sums the total weight of all paths in a graph where the paths are written in a very direct style as a list of nodes visited along the path.

```
568 goal += start(X_1) * \beta([X_1|X_S]).

569 \beta([X_1,X_2|X_S]) += edge(X_1,X_2) * \beta([X_2|X_S]).

570 \beta([X]) += stop(X).
```

We apply the right-corner transformation to this program with σ equal to all the rules. We get the following after useless-rule elimination:

```
571 (Z / \beta([X_1|X_S])) += start(X_1) * Z / goal.

572 (Z / \beta([X_2|X_S])) += edge(X_1, X_2) * Z / \beta([X_1, X_2|X_S]).

573 (Z / stop(X)) += Z / \beta([X]).

574 (Z / Z).

575 goal += stop(X_1) * goal / stop(X_1).
```

After specialization:

```
576 (goal / path([X<sub>1</sub>|Xs])) += start(X<sub>1</sub>) * goal / goal.
577 (goal / path([X<sub>2</sub>|Xs])) += edge(X<sub>1</sub>, X<sub>2</sub>) * goal / path([X<sub>1</sub>, X<sub>2</sub>|Xs]).
578 (goal / stop(X)) += goal / path([X]).
579 (goal / goal).
580 goal += stop(X<sub>1</sub>) * goal / stop(X<sub>1</sub>).
```

When we run this program, the tail of the list Xs in goal/path([$X_2|Xs$]) is never instantiated. Thus, we can abbreviate goal/path([X|Xs]) \mapsto starpath(X), and we recover the efficient algorithm from Example 27.

7.6 Discussion

Implementation. We have implemented and tested these transformations. They are available in our system: https://github.com/timvieira/dyna-pi. Using these transformations in the context of an automated optimizer (§8) has stress-tested them since buggy transformations often lead to fast programs (e.g., the empty program), and search probes the possible arguments to the transformations.

Eisner and Blatz (2007). The original paper on transforming semiring-weighted logic programs is Eisner and Blatz (2007). They extended the fold/unfold

transformations (Tamaki and Sato, 1984), and the magic sets transformations (Ramakrishnan, 1991) from the unweighted setting to the weighted setting and proposed a novel speculation transformation.

Single vs. Multiple Semirings. Eisner and Blatz (2007) allowed programs to use more than one semiring to combine values. They supported this case by transforming semiring program fragments; our work should also be compatible with semiring fragments. As we discussed in §4, our semantics are based on sums over derivations, and Eisner and Blatz (2007)'s are based on the weaker notion of fixpoints. The sum-over-derivations semantics is deeply tied to the distributive property of the semiring. Mixing semirings means that property is not always available. 195

Eisner and Blatz (2007)'s extensions of the fold and unfold transformations to the weighted setting were limited to *basic* folds. Consequently, their exploration of what fold and unfold can accomplish was severely limited. In particular, they never discovered self-supporting recursive definitions like startpath (Example 27). We generalized Eisner and Blatz (2007)'s treatment of fold/unfold to allow fold and unfold with respect to auxiliary programs, such as prior versions of the program in a sequence of transformations. The generalization to auxiliary definitions was previously observed to be powerful in the unweighted setting (e.g., Tamaki and Sato (1984), Roychoudhury et al. (1999), Bossi et al. (1992), and Pettorossi

¹⁹⁵A notable exception is when the mixed-semiring program is *stratified*—meaning that we can partition its non-<u>0</u> items into strongly connected components (as we did in §5.4) such that each component uses a single semiring.

and Proietti (1994)). Examples like Example 27 and Example 28 are impossible with only basic fold. We give examples of transformed programs going from non-terminating to terminating and exponential to polynomial time. Unlike the limited examples presented in Eisner and Blatz (2007) (and Gildea (2011)), which only employed simple, superficial "hook tricks," our approach led to the discovery of self-supporting recursive definitions. However, some additional care must be taken to prevent ill-founded circular definitions. We characterized the ill-founded definition problem in §7.4.2.2 as the nontermination of a derivation mapping. We proposed a practical system for avoiding ill-founded definitions in §7.4.2.4 that builds on Roychoudhury et al. (1999).

⇒ **Speculation as a Hoisting Transformation.** We deepen the discussion of the speculation transformation by characterizing it as a structural change to the derivations of a program. We also connected it to the better-known left-corner grammar transformations (Rosenkrantz and Lewis, 1970; Johnson, 1998; Johnson and Roark, 2000), which we generalized to Dyna programs (§7.5.2).

Correctness Proofs. Eisner and Blatz (2007) never formally proved the correctness of (basic) folding or unfolding; we have filled this gap in the semiring-weighted logic programming literature.

Contributions to the Fold and Unfold Literature. We analyze the fold—unfold transformation's ability to preserve the entire set of derivations for each item. Prior work¹⁹⁶ focused on the preservation of *at least one* derivation

¹⁹⁶E.g., Tamaki and Sato (1984), Maher (1987), Kawamura and Kanamori (1990), Bossi et al. (1992), Roychoudhury et al. (1999), and Pettorossi and Proietti (1994).

for each item, which is insufficient for our semiring-weighted setting.

Extending Roychoudhury et al. (1999). Our practical fold/unfold safety test (§7.4.2.4) builds directly on Roychoudhury et al. (1999) for unweighted logic programming. We proved that their technique extends directly to weighted logic programming. However, we sought a different notion of correctness: strong equivalence. As a result, we provide an alternative lens on why their technique is correct. We have also structured our presentation (and proofs) to disentangle the structure of the transformation from the static analysis technique used to check its safety. Additionally, §7.4.2.4 makes their approach's static analysis modeling assumptions more explicit than in other work. One minor contribution is that we have extended their method to support the generalized unfold transformation, which was straightforward but missing in their work. Lastly, as a practical matter, we proposed using constraint solvers to eliminate the need to specify an initial measure, which improves our ability to approve certain transformations that would otherwise be rejected.

A Notable Omission: The Magic Sets Transformation. Eisner and Blatz (2007) gave a weighted generalization of the magic sets transformation (Ramakrishnan, 1991), a well-known transformation to the Datalog community as it is a simple method for partially specializing a program to some specific queries. An exciting application of the magic sets transformation is that it can recover the top-down filtering part of Earley's algorithm when applied to CKY (Minnen, 1996; Eisner and Blatz, 2007). Unfortunately, we do not support Eisner and Blatz (2007)'s version of magic sets because the transformed program mixes

semirings. Support for mixed semiring programs is interesting future work. However, a different perspective on the magic sets transformation is that it is a recipe for prioritization of agenda-based evaluation ($\S 5.1.3$). It is possible to derive tighter heuristics than the simple boolean heuristic given by magic sets. Felzenszwalb and McAllester (2007) and Prieditis (1993) show how to use the outside context of a *relaxed* program to prioritize the fine-grained inside pass. The relaxed programs can be automatically discovered and optimized by searching over which constraints and variables to drop. The values that it computes for prioritization can be more expressive than booleans. Moreover, in the case of superior semirings, we can even automatically derive admissible heuristics for A^* -like decoding. In future work, we would like to explore this direction.

Summary

We presented several program transformations that can be used to optimize a given program. We have many examples where the transformations dramatically improve the running time. In §8, we will use search to automatically find a transformation sequence that minimizes the program's running time within a search budget.

Chapter 8

Searching For Faster Dynamic

Programs

This chapter develops a tool for automating the discovery of faster programs. At a high level, our approach is straightforward: given an initial, correct program, we will search for a sequence of meaning-preserving, source-to-source transformations (§7) that minimizes a given cost function, such as the worst-case runtime bounds of §6.3 or the average execution time of the program on a given input workload. We will discuss the benefits and drawbacks of several cost functions. We will describe how we structure our transformation sequence search space. We conclude with experiments demonstrating that our proof-of-concept system is capable of replicating the examples given throughout the thesis as well as many of the missed speed-ups found in the literature (outlined in Example 1).

8.1 The Graph Search Problem

We consider an abstract **graph search problem**, $\langle \mathcal{S}, s_0, \mathcal{T}, \text{transitions}, \text{cost} \rangle$, where \mathcal{S} is a state space, transitions : $\mathcal{S} \to 2^{\mathcal{S}}$ is a transition function, $s_0 \in \mathcal{S}$ is an initial state, $\mathcal{T} \subseteq \mathcal{S}$ is a set of terminal states, and $\text{cost}: \mathcal{T} \to \mathbb{R}_{\geq 0}$ is a cost function on the terminal states. The search algorithm will treat the cost and transitions functions as black boxes. The goal of the search problem is to find the (reachable) terminal state in the graph that has the lowest cost, $s^* = \min_{s \in \mathcal{T}} \text{cost}(s)$.

Our program optimization problem (depicted in Figure 8.1) can be straightforwardly mapped into this notation. Our states $s \in \mathcal{S}$ are programs. The initial state s_0 is the initial program. Each transition is the application of a semantics-preserving program transformation (§ 7). This means every state reachable from s_0 in this graph is a semantically equivalent program. In our setting, every state is terminal. In §8.1.1, we will describe the search algorithm that we use to traverse this graph. In §8.1.2, we will discuss the details of how we structured the search space used in our experiments (§8.3). In §8.2, we will discuss the cost function.

¹⁹⁷In principle, the cost function does not have to be real-valued; it can be any mapping from states to a totally ordered set.

¹⁹⁸To use the generalized fold/unfold transformations, the state must be augmented with the metadata required to check the measure-based safety conditions of §7.4.2.4.

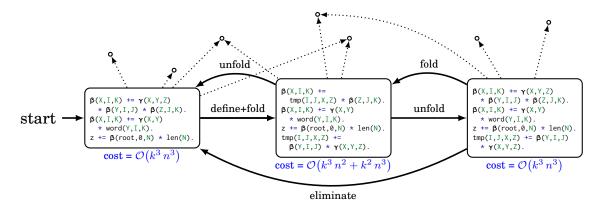


Figure 8.1: Depiction of the program optimization graph search problem. The program used in this figure is CKY (Example 6). Nodes are Dyna programs (§ 2). The node pointed to by "start" indicates the user's program. Edges are program transformations (§ 7). Costs are derived by program analysis (§ 6). Only a tiny subset of the nodes and edges in the search graph is shown. The dotted unlabeled outgoing edges represent additional transformations we did not elaborate on in the diagram to reduce clutter.

8.1.1 Prioritized Search Algorithm

Our search algorithm is based on best-first search, which is often used to solve difficult combinatorial optimization problems (Russell and Norvig, 2020). Pseudocode is given in Algorithm 30. The basic idea of the best-first search algorithm is to maintain a set of candidate states $\mathcal Q$ to explore next and to expand the state that is estimated to be the most promising based on a search policy π that selects among the candidates. Once a candidate is selected, its successor states are enumerated according to transitions and added to the candidate set $\mathcal Q$ (if they have not previously been added).

Our overall search algorithm search_{π} is a kind of *anytime* algorithm (Zilberstein, 1996) that traverses the graph until the search budget is exceeded. While search_{π} traverses the graph, it keeps track of the least cost terminal state it has found so far. If the search budget is reached before the algorithm has explored

Algorithm 30 search: Algorithm for a given graph search problem $\langle \mathcal{S}, \mathcal{A}, s_0, \mathcal{T}, \text{transitions}, \text{cost} \rangle$. Requires a search policy π to select which candidate state to expand. We also require that a search budget is specified, such as an upper limit on the total execution time spent searching.

```
1. def search_{\pi}(s_0):
      \mathcal{Q} \leftarrow \{s_0\}
                                                                > Set of candidate states to explore next
      visited \leftarrow \{s_0\}
                                                          ▷ Set of states that have already been visited
      while |\mathcal{Q}| > 0 and the search budget has not been exceeded:
         s \leftarrow \pi(\mathcal{Q})
                                                                 ▶ Pick any state among the candidates
         Q.remove(s)
         if s \in \mathcal{T}:
            evaluate cost(s)
         for s' \leftarrow \text{transitions}(s):
            if s' \notin \text{visited}: \triangleright Add s' as a candidate, if we have not already done so in the past
10
                Q.add(s')
               visited.add(s')
     return the lowest cost state evaluated
```

the entire graph, it will return the best solution found up until that point. This makes it an *anytime* algorithm because it can provide a solution even if it has not fully explored the search space (i.e., before |Q| = 0). The search policy π governs the order in which the states of the graph are explored. Abstractly, π takes the set of candidate states $Q \subseteq S$ and returns an element from Q, possibly randomly. It can encode a variety of search orders (e.g., depth-first search, breadth-first search, greedy search, and A^*). In practice, we use a priority queue to implement Q and π , and we compute the priority of a given state when we add the state to Q (line 30.11).

This algorithm can guarantee that the best reachable state is eventually found if the search budget is sufficiently large and π does not *starve* any candidate state.¹⁹⁹ A simple policy that guarantees this is the breadth-first policy,

¹⁹⁹ Starving means that the state says on the priority queue indefinitely.

which enumerates states in first-in-first-out order (FIFO); or, equivalently, by the length of the path by which it was first discovered. This strategy does not have good anytime performance. On the other hand, the greedy strategy of always selecting the cheapest state often has good anytime performance, but it may starve states. Therefore, we use a hybrid strategy that interleaves greedy and FIFO choices to get the benefits of both.

8.1.2 Search Space Details

The most basic version of the search graph would say that all valid transformations from §7 are available at all times.²⁰⁰ However, this would result in an astronomical number of states to explore, and it would ignore some useful problem structure. This section provides the details of how our search space is structured.

Eureka definitions. For proposing Eureka definitions (§7.3), our automated system considers variable-elimination actions. To enumerate the available variable-elimination definitions, we consider each rule one at a time. For a given rule R, the variable $v \in \text{vars}(\mathbf{R})$ can be usefully eliminated if it does not appear in all of the subgoals in the rule's body and it does not appear in the head of the rule. Thus, to enumerate all variable-elimination definitions, we

 $^{^{200}}$ In the sections of § 7 describing each transformation, we discussed the conditions for the transformations to be valid.

²⁰¹In principle, other Eureka definitions §7.3 be straightforwardly added within our search framework as additional transformations.

²⁰²More formally, the set of variables that can be eliminated is $\operatorname{elim}(\mathbf{R}) \stackrel{\text{def}}{=} \{v \mid \operatorname{factors}(\mathbf{R},v) \neq \operatorname{body}(\mathbf{R}), v \notin \operatorname{head}(\mathbf{R})\}$ where $\operatorname{factors}(\mathbf{R},v) \stackrel{\text{def}}{=} \{y \mid y \in \operatorname{body}(\mathbf{R}), v \in \operatorname{vars}(y)\}$. If any rule \mathbf{R} of the program $\mathfrak p$ contains variables that can be eliminated ($|\operatorname{elim}(\mathbf{R})| > 0$), then eliminating any variable $v \in \operatorname{elim}(\mathbf{R})$ by folding factors(\mathbf{R},v) out of \mathbf{R} into a new rule reduces that \mathbf{R} 's runtime bound, which may reduce (and never increases) the runtime bound of the program. When

identify each variable that can be eliminated and propose a new definition using the recipe detailed in §7.3.

Unfold. Recall that the generalized unfold transformation (Definition 35) is parameterized by a source program $\mathfrak{p}^{(s)}$, rule index $r \in \mathfrak{p}^{(s)}$, a subgoal index of r, and auxiliary definitions $\mathfrak{p}^{(d)}$. Under this parameterization, it is straightforward to enumerate all instances of unfold: any rule and any non-input subgoal in that rule are candidates for unfolding. The choice of $\mathfrak{p}^{(d)}$ is discussed later. Generalized unfolds must pass the measure-based safety test.

Fold. Recall that the generalized fold transformation (Definition 36) is parameterized by a source program $\mathfrak{p}^{(s)}$, auxiliary definitions $\mathfrak{p}^{(d)}$, a new rule r, and a set of rules $S \subseteq \mathfrak{p}^{(s)}$ that will be replaced.

- We generate candidate rules r as follows: for each pair $\langle \mathbf{S}, \mathbf{D} \rangle \in \mathfrak{p}^{(s)} \times \mathfrak{p}^{(d)}$, we attempt to replace subproducts of S's body by \mathbf{D} .²⁰³ From there, we obtain a set of r candidates and filter out those that do not satisfy the reversibility condition of the fold transformation.
- The set S is determined by r and the position j that we unfold in r's body in the generalized reversibility test. In the generalized reversibility test, we have a generalized unfold transformation. That unfold transformation adds a set of rules, and that set of rules is S.

more than one variable can be eliminated ($|elim(\mathbf{R})| > 1$), the order in which the variables are eliminated will affect the eventual runtime bound. Unfortunately, finding an optimal sequence of variable-elimination steps is NP-hard by reduction from variable-elimination ordering in probabilistic graphical models (Eisner and Blatz, 2007; Gildea, 2011).

²⁰³The relevant algorithms for enumerating subproducts are given in §C.2. Once the rule bodies are aligned, we replace the matched subproduct with the head of D. We will subsequently check the validity of this replacement using the reversibility condition.

• Generalized folds must pass the measure-based safety test.

Auxiliary definitions used for generalized fold/unfold. For efficiency, we limit the auxiliary definitions considered by generalized folds and unfold transformations to the original definitions of each item.²⁰⁴ We have found our choice sufficient for our benchmark problems.

Reducing the number of safety checks. We found it more efficient overall to check the fold/unfold safety conditions after the transformed program is selected by π . This results in significantly fewer calls to the constraint solver. This way, we only have to check the safety conditions of the *most promising* transformations rather than *every* transformation. Technically, this means that a state is only terminal (i.e., a member of \mathcal{T}) after it passes the safety test. Also, we do not continue transforming states deemed unsafe by the test.

Hoisting transformations. Recall the hoisting transformations (i.e., speculation and left-corner) are parameterized by a source program $\mathfrak{p}^{(s)}$ and a subset of rules $\sigma \subseteq \mathfrak{p}^{(s)}$. For speculation, we have an additional parameter $z \in \mathbb{H}$ that describes the type of the item. In the non-commutative case, the choice of left vs. right is also a parameter. In the commutative case, we consider moving any of the subgoals in each of the rules in σ to be the leftmost subgoal as an additional parameter. Unfortunately, the hoisting transformations are difficult to search over because there are far too many of them available: the number of valid hoisting transformations scales exponentially with the number of rules—and it

²⁰⁴In principle, any program visited in the entire search space can be used as an auxiliary definition. However, that yields far too many options, and these options do not appear to increase the set of states that are reachable in practice.

is even worse if the program's semiring is commutative! This may be tolerable for a small set of programs, each with only a few rules, but incorporating it into search will result in far too many candidate states. We have included methods for enumerating these transformations in the software library, but they must be used sparingly, or else the search space grows too large.

Program specialization. Recall that the abbreviation transformation (§6.2) can speedup a given program, but may do so by eliminating useless items (i.e., making their values 0). However, our transformation framework (§7) assumes all programs are (strongly) equivalent; thus, we may only apply abbreviation as a post-processing step. The cost function, on the other hand, may wish to account for the cost of the abbreviated program. However, running abbreviation often may slow down search.

Rule to-do list. The unfold and Eureka generation methods we consider are centered around a specific rule in the current version of the transformed program. Applying transformations to rules $R, R' \in \mathfrak{p}^{(s)}$ in either order will get the same result if neither transform makes the other one impossible. Thus, we consider transformations in a canonical order. Each state will now consist of a program together with a **to-do list** of rules that can still be transformed. The possible actions at that state consist of removing the top rule R from the list (declining to transform R) or applying the next transformation. Applying such a transform may delete and/or add program rules, which are correspondingly deleted from the list and/or added at the bottom. (If the list is empty, no more actions are possible.) This design reduces the branching factor by an amount

proportional to the number of rules. Of course, a potential downside is that it makes the search tree deeper by a factor proportional to the number of rules.

8.2 Cost Functions

There are two broad families of cost functions: empirical and analytical. Our experiments will only use analytical bounds; we will explain our position here.

Empirical cost functions. The most direct cost function to consider optimizing is the empirical execution time of running a given program on a specific workload of representative inputs (e.g., running a transformed parser on actual sentences from the Penn Treebank (Marcus et al., 1993)) with a specific solver configuration (e.g., prioritized fixpoint iteration with SCC-based prioritization) on a specific computer (e.g., your laptop).

Analytical cost functions.²⁰⁵ Analytical runtime bounds, such as worst-case runtime bounds, are the method of choice for the theoretical comparison of algorithms in computer science. Indeed it is how we have compared algorithms throughout this dissertation. In our setting, worst-case, analytical running time bounds can be derived by the approach that we gave in §6.3, which automates McAllester (2002)'s well-known meta-theorem.²⁰⁶ We will additionally provide a

²⁰⁵Although we do not explore it in this work, we briefly mention that it is possible to derive average-case runtime estimates from statistical models of the input data. Such models are used in database query optimizers (Deshpande et al., 2007; Chekuri and Rajaraman, 1997; Leis et al., 2017).

²⁰⁶This theorem is used throughout natural language processing to analyze algorithms, e.g., Gildea (2011), Nederhof and Satta (2011), Gilroy et al. (2017), Melamed (2003), Kuhlmann (2013), Nederhof and Sánchez-Sáez (2011), Büchse et al. (2011), Lopez (2009), and Eisner and Blatz (2007).

cheaper-to-compute proxy later in this section.

Direct execution can be expensive. Empirical cost functions can be very expensive to evaluate—especially when the initial program is inefficient. For example, an initial program could have an infinite running time, but its optimized version might be finite. Similarly, evaluating the execution time of an $\mathcal{O}(n^{1000})$ program, unsurprisingly, requires $\mathcal{O}(n^{1000})$ time. To mitigate the issue of some programs running very long on certain inputs, it is common to impose a time limit on the execution of any given program. We also have the challenge that the runtime measurements must be averaged over many runs to estimate the expected running time precisely. Often a deterministic proxy measurement is preferable. In our setting, the observed number of total prefix firings or the program's rules (§6.3) strongly correlates with the execution time of fixpoint iteration algorithms.

Overfitting the workload. The direct execution approach requires us to specify particular inputs to the program; this has the downside in that we can *overfit* to these inputs. We may also overfit to specific details of how the program was executed (e.g., hardware-specific details). Optimizing analytical measures, or other proxies for execution time, can lead to *underfitting*.

Practical efficiency vs. lack of fidelity. Practical efficiency can easily outweigh the lack of fidelity. Consider evaluating the cost of an $\mathcal{O}(n^{1000})$ program: evaluating the wallclock time takes $\mathcal{O}(n^{1000})$ time, whereas assessing the runtime complexity does not depend on n in this way. Optimizing the analytical runtime can provide a crucial speedup over direct execution, enabling a more

exhaustive search in practice. Additionally, it sidesteps the need to specify a specific workload. Of course, the efficacy of comparing analytical bounds is highly dependent on the fidelity of the analysis used to derive the bound: comparing very loose bounds can lead to incorrect conclusions about which algorithm is more efficient.²⁰⁷ The assumptions behind the bounds are also important, e.g., the tightness of the input-type specification.

Solver configuration. We are considering the optimization of programs, not solver configurations. However, we want to highlight that the solver configuration, including the choice of the solver itself, is an important choice when developing an efficient, end-to-end algorithm for a given problem.²⁰⁸

• The choice of solver from §5 is important: prioritized fixpoint iteration and Newton's method are good choices, but they have their pros/cons that depend on the program and its input data. Additionally, each of the solvers has its configuration options. For example, the prioritized fixpoint iteration algo-

 $[\]overline{g}^{207}$ We caution that comparing upper bounds does not lead to sound conclusions about which system is better. In principle, given $f_{lo} \leq f \leq f_{hi}$ and $g_{lo} \leq g \leq g_{hi}$ we can only correctly conclude that f is definitely smaller than g iff $f_{hi} \leq g_{lo}$. In other words, the correct conclusion requires both an upper and lower bound. Unfortunately, we do not have a method for deriving meaningful lower bounds on running time. Nonetheless, the comparison can be accurate if the runtime bounds are tight and characterize the typical behavior.

²⁰⁸Many high-performance computing systems are tuned by optimizing execution time. Some well-known examples are numerical libraries, such as ATLAS (Whaley, Petitet, et al., 2001), PhiPAC (Bilmes, Asanović, et al., 1996), FFTW (Frigo and Johnson, 1998), and SPIRAL (http://www.spiral.net), where execution time is used to adapt low-level implementation details to specific hardware. These details include loop {order, tiling, unrolling} and whether to use instruction-level parallelism. One of the primary reasons these libraries are so popular is that performance improvements on one computer do not necessarily generalize to other architectures or all workloads; that is why they tune the implementation for specific hardware. A portable performance research community has formed around tuning software to be as fast as possible on a given computing platform. We also mention the configuration of solvers for computationally hard problems. Hutter (2009) used black-box optimization techniques to tune the parameters of solver heuristics for mixed-integer programming solvers and satisfiability solvers to optimize efficiency for specific problem workloads. This approach won several competitions, demonstrating the effectiveness of tuning solvers to a specific workload.

rithm's execution time will depend on the choice of the prioritization function. The choice of convergence threshold (if applicable) affects the running time for cyclical programs using the SCC-based approach. Newton's method has overhead that fixpoint iteration does not have, but it is quicker at solving cyclical programs.

• For any given program and solver, many low-level details affect the execution time. For example, memory layouts (e.g., row-order or column-order layout of a dense array in memory), sparse vs. dense representations of relations (e.g., hash tables vs. arrays), and indexes on relations (including sorted order) can have a dramatic effect on the running time in practice. However, they will not manifest in the analytical runtime bounds.

Fortunately, improvements to the program are generally orthogonal and complementary to the lower-level solver configurations. Our perspective is that we should get the high-level algorithm right—i.e., find an algorithm with a good big- \mathcal{O} running time—before optimizing the lower-level choices. From there, we can investigate low-level tricks, such as memory layouts, that improve the constant factors and average execution time.

8.2.1 Fine-grained Runtime Bounds

In §6.3, we devised an automated method for assessing a given program's worst-case time complexity as a function of user-defined input-size parameters.²⁰⁹

²⁰⁹The validity of this bound was discussed in §6.3, but we remind the reader that the bound assumes that the program is acyclic for the bound to be valid. Additionally, the bound assumes

Comparing bounds. The analytical bounds provided by the method need to be compared. The most obvious approach is to evaluate the runtime expressions on specific values for the input-size parameters. For runtime bounds in a single variable, it is possible to determine the faster of two runtime bounds in the sense of big- \mathcal{O} using a computer algebra system, such as SymPy (Meurer et al., 2017) or Mathematica (Wolfram, 2003). Comparing multi-variable runtime expressions in this way is trickier. One option is for the user to define the input-size parameters as a function of a single parameter t (e.g., $k = \log(t)$, n = t), and reduce to the single variable case.

Efficiency. The method from §6.3 can be slow for a few reasons:

- For some programs, abstract forward chaining may need to do a lot of work. At one extreme, a program could have no input data; thus, abstract interpretation will essentially solve the relaxed, booleanized version of the program. Thus, similar to concrete execution, abstract forward chaining may also require limits on its execution, as discussed in (§6.1.2.5).
- Overhead: Constraint propagation systems can be slow if the constraint propagation rules are very productive (produce a lot of consequences). Additionally, the method from §6.3 optimizes subgoal-ordering decisions, which can be inefficient for rules with long bodies.²¹¹

²¹¹A suboptimal subgoal ordering could be used to speed things up.

that prioritized fixpoint iteration is used with SCC-based prioritization (§5.1.4).

²¹⁰These computer algebra systems can compare the asymptotic behavior of two functions f(n) and g(n) in the sense of big- \mathcal{O} : $f \leq g \iff f = \mathcal{O}(g) \iff \limsup_{n \to \infty} \frac{f(n)}{g(n)} < \infty$. If the runtime bound is a simple polynomial, this is just a comparison of the largest degree terms. However, the runtime bounds of §6.3 may include \min in addition to summation and multiplication, which makes the comparison more detailed; hence we suggest the computer algebra system.

These inefficiencies motivate the coarser-grained analytical bound that we describe in §8.2.2.

8.2.2 Coarse-grained Runtime Bounds

In this section, we will describe a very fast cost function based on a folk theorem from the Datalog community with a long history of use. The folk theorem says that, under certain conditions (discussed later), the running time of forward chaining execution of a given program is, at worst, linear in the number of ways to instantiate its rules, i.e., bind the variables to constants. A relatively simple bound on rule instantiations is available if we can establish that each variable in the program can be bound in at most η different ways. Given some other conditions discussed at the end of this section, the number of ways to instantiate a rule with k variables is bounded by $\mathcal{O}(\eta^k)$. Program \mathfrak{p} 's total runtime is $\mathcal{O}(\eta^{\text{degree}(\mathfrak{p})})$ where $\text{degree}(\mathfrak{p})$ is the maximum number of variables in any rule of \mathfrak{p} . We can, therefore, use degree as a cost function to minimize during search, and it is very efficient to compute and easy to implement.

As an example, consider CKY. The degree bound says that the program below runs in $\mathcal{O}(\eta^6)$, as the first rule must sum over 6 variables (X, Y, Z, I, J, K).

```
600 \beta(X,I,K) += \gamma(X,Y,Z) * \beta(Y,I,J) * \beta(Z,J,K).

601 \beta(X,I,K) += \gamma(X,Y) * \beta(Y,I,J).

602 \beta(X,I,K) += \gamma(X,Y) * word(Y,I,K).
```

This is a coarse-grained analysis: the runtime is usually given more specifically as $\mathcal{O}(n^3\,k^3)$ where n is the number of sentence positions, and k is the number of $\overline{}^{212}$ This is not to be confused with the number of non- $\underline{0}$ instantiations of the rule.

grammar symbols. (This finer-grained bound can be achieved using type analysis §6.3: the intuition is that the variables I, J, K can each be bound in $\mathcal{O}(n)$ ways while X, Y, Z can each be bound in $\mathcal{O}(k)$ ways.) However, the simpler analysis $\mathcal{O}(\eta^6)$ gives us a single exponent to reduce, namely 6. Below is a sped-up variant, which sums over the variable Y separately from K:

```
603 \beta(X,I,K) += tmp(I,J,X,Z) * \beta(Z,J,K).

604 tmp(I,J,X,Z) += \gamma(X,Y,Z) * \beta(Y,I,J).

605 \beta(X,I,K) += \gamma(X,Y) * \beta(Y,I,J).

606 \beta(X,I,K) += \gamma(X,Y) * word(Y,I,K).
```

This variant is more efficient as its running time is $\mathcal{O}(\eta^5)$. It is also more efficient under the finer-grained analysis, $\mathcal{O}(k^2 n^3 + k^3 n^2)$.

For programs with multiple rules, we can refine the degree analysis to measure the degree of each rule in the program rather than summarize it by the maximum degree. Then, we can say that a program is faster if its list of rule degrees—sorted in descending order—is lexicographically smaller than that of the other program. This has the benefit of reducing lower-order terms in the running time, not just the highest-order term.

Correctness of the bound. The degree analysis of a Dyna program only leads to a *valid* \mathcal{O} -expression under some conditions, which we will now discuss.

- 1. The degree bound assumes the *grounded* program to be acyclic. Cycles slow down forward chaining because we must iterate to a fixpoint (§5). Generally, the number of iterations required to reach a fixpoint is data-dependent.
- The degree bound assumes that all of the relations in the program are finite.This prevents the user from encoding infinite sets, such as the Peano integers.

Although technically, the bound is not *wrong*; it is just that η is infinite.

3. The degree bound also assumes constant-time semiring operations.

Recall from §6.3 that the fine-grained runtime bound makes assumptions (1) and (3). However, the fine-grained bound is generally tighter as it can account for certain kinds of sparsity, and its bounds do not assume all variables have the same domain size. The fine-grained system can *check* if (2) is true by type inference, so it is not *assumed*.

Looseness. The degree analysis might be loose for many reasons. The upper bounds derived using our methodology assume that relations are dense and that all variables have a single type. We note that the degree analysis becomes tighter after performing program specialization (§6.2). For example, if two variables in the program are always equal, we should not count them as separate, but the degree bound would. Fortunately, specialization can equate these variables, which improves the degree bound. Similarly, if a variable is always bound to a finite set of values (or it is always uninstantiated), it should not count at all in the degree; specialization transforms the program not to have such variables. Despite being a crude approximation, minimizing the degree of a program often works well. We will see in §8.3 that simply optimizing degree is sufficient to recover several asymptotic speedups noted in the NLP literature (see § 1).

²¹³To speed up search, our experiments only run specialization on the final program.

8.3 Experiments

We aim to devise a system for automatically improving typical dynamic programming problems—especially those encountered in natural language processing applications. To evaluate whether we achieved this goal, we devised a set of benchmark programs to see how well our proposed approach works. Our tests include most of the missed speedups that we mention in §1 as well as many of the examples that we described throughout the dissertation; the specific set of programs is provided in §E.

Experiment details.

- Results are summaried in Table 8.1.
- The name of each benchmark is hyperlinked to a section in §E that contains further information about the benchmark problem, including the Dyna code of the initial program and any notes about where the program appears in the literature. For each of the benchmark problems, we have manually specified its optimal degree. We assume a commutative semiring in our experiments.
- We show for each benchmark the initial, found, and optimal degrees as well as search time: the time (in seconds) that it took search to find the optimal solution (if one was found). Cells marked with "—" indicate that the optimal-degree solution was not found within the search-time budget of 1 minute.
- A check mark () signifies that the benchmark problem was solved to optimal-

ity with respect to the degree cost function.

• For the search policy π , we use a hybrid greedy–FIFO strategy (suggested in §8.1.1) where we alternate between greedy and FIFO actions.²¹⁴ The greedy policy is prioritized by degree.

Discussion.

- These tests are by no means exhaustive, but they show that our proof-of-concept system can automate a variety of speedups. We see that in all but two cases, search is able to find the optimal degree program within a few seconds. We believe that we could significantly improve this already quick time by implementing each of the program transformations more efficiently. For example, in the implementation used for these experiments, there is a lot of unnecessary copying, which can be expensive as programs become larger. Improving the efficiency of the program transformations would improve our ability to search the state space more quickly and exhaustively.
- The benchmarks {path-list, even-odd-peano, explicit-pda} are unlike the others in that each of their initial versions has an infinite-size variable domain (i.e., $\eta = \infty$). However, they can each be optimized so that they have finite domains. Despite the degree bound being unaware of the change in η , the programs found by search for path-list and even-odd-peano are optimal.

²¹⁴We did not optimize the ratio of greedy and FIFO actions since it could lead to overfitting this small set of benchmark problems. We did find that the purely FIFO policy increased average search time significantly.

name	initial degree	found degree	optimal degree	search time (sec)
arc-eager	6	✓ 3	3	1.375
bad-chain-05	3	✓ 2	2	0.398
bad-chain-10	3	✓ 2	2	1.289
bar-hillel	10	✓ 8	8	0.013
bilexical-labeled	8	√ 7	7	0.014
bilexical-unlabeled	5	✓ 4	4	0.049
chain-05	6	✓ 2	2	0.016
chain-10	11	✓ 2	2	0.090
chain-expect	3	✓ 3	3	0.003
cky+grammar	6	✓ 3	3	1.424
cky3	6	✓ 5	5	0.003
cky4	8	✓ 6	6	0.012
edit	6	✓ 4	4	0.008
even-odd-peano	1	~ 0	0	5.812
explicit-pda	8	6	5	_
hmm	5	✓ 4	4	0.003
itg	9	✓ 8	8	0.011
path-list	3	✓ 2	2	1.164
path-start	3	✓ 2	2	8.707
semi-markov	4	✓ 3	3	0.002
split-head-EB	5	4	3	_
split-head-J	4	✓ 3	3	0.780
tensor-decomp-parser	6	✓ 4	4	1.916

 $\textbf{Table 8.1:} \ Experimental \ demonstration \ of our search-based \ program \ optimizer \\ on \ the \ benchmarks \ described \ in \ \S E.$

- We also highlight that our system was unable to solve split-head-EB benchmark to optimality. We did manage to speed it up $(\mathcal{O}(\eta^5) \to \mathcal{O}(\eta^4))$. However, we are unable to achieve the optimal $\mathcal{O}(\eta^3)$ algorithm because we do not have hoisting transformations in our search space (discussed in §8.1.2). Unfortunately, we have not yet found an efficient strategy for searching over hoisting transformations. However, we have verified that manually initializing the transformation search with one round of speculation transformations on the initial program is sufficient for our algorithm to recover the optimal degree algorithm. The issue is that making these transformations available in other states causes the search space to grow too much.
- The explicit-pda benchmark appears to be out of reach for our space of program transformations. We describe how to optimize it in §E.12. The missing trick appears to be a generalization of speculation.

Summary

In this chapter, we developed a tool for automatically discovering faster programs. We do so by searching for a transformation sequence that minimizes a given cost function, which we cast as a kind of graph search. We gave a simple anytime algorithm for performing program optimization on this graph. We discussed the benefits and drawbacks of using various cost functions for our optimization problem. We concluded with an experimental demonstration of how our prototype system can optimize several benchmark problems drawn from the literature on overlooked speedups (Example 1). In all but two of our benchmarks, search finds the fast algorithm in a matter of seconds. However, there is room for improvement. A particular issue we ran into experimentally was the difficulty of searching over hoisting transformations because too many of them were available. We also showed an example that is out of reach for our current space of transformations—suggesting that future work will want to broaden the available transformations. Nonetheless, we are pleased to see that so many of the speedups found in the literature can be rapidly found automatically. Thus, we believe this tool already provides significant value as these previously overlooked speedups and similar ones will no longer be missed when using our tool.

Chapter 9

Related and Future Work

9.1 Related Work

Dyna 1. We have drawn heavily on Eisner, Goldlust, et al. (2005) and Eisner and Blatz (2007). We have discussed these connections extensively throughout this dissertation.

Dyna 2. The Dyna 2 language (Eisner and Filardo, 2011) is closely related to the Dyna 1 language that we have discussed in this work, as Dyna 1 programs are contained in the set of Dyna 2 programs. We have provided a short section in §D describing the language and possible extensions of our work.

Algorithm = Logic + Control. The decomposition of an algorithm (for parsing or otherwise) into a logical specification and controller is due to Kowalski (1979), who wrote that "Algorithm = Logic + Control." Later, in the context of parsing algorithms, Shieber et al. (1995) wrote:

"Although the deduction systems do not specify detailed control structure, the control information needed to turn them into full-fledged parsers is uniform, and can therefore be given by a single deduction engine that performs sound and complete bottom-up interpretation of the rules of inference."

The logical specification of algorithms was further studied in many works, such as McAllester (2002), Ganzinger and McAllester (2002). However, it was truly espoused by the natural language parsing community after the seminal papers of Pereira and Warren (1983) and Shieber et al. (1995). Outside the NLP community, the programming languages community widely uses logical deduction notation to specify operational semantics (Plotkin, 2004) and program analyses (Whaley, Avots, et al., 2005; Ullman, 1988; Dawson et al., 1996; Reps, 1993; Esparza et al., 2010).

Parsing as Deduction. Pereira and Warren (1983) wrote a seminal paper to the NLP and parsing community, which cast parsing as logical deduction. Later, Shieber et al. (1995) established deduction systems as a valuable framework for concisely specifying parsing algorithms.

Goodman (1998) modernized the deduction framework to accommodate what was at that time a budding interest in *statistical* parsing. He showed that many quantities of interest, such as the total probability, the most-likely parse, and the set of most-likely parses, could be computed using the same weighted inference algorithm where differences are encapsulated into a choice of semiring

weights.²¹⁵ This was an important discovery, and it aligns with work in graphical models that was going on around the same time (Aji and McEliece, 2006). At this point, the need for a programming language for specifying *weighted* deductive system was evident. Thus, the Dyna notation emerged (Eisner, Goldlust, et al., 2005; Eisner, Goldlust, et al., 2004).

Deduction systems have been used to express several formalisms:

- Earley parsing (Pereira and Warren, 1983)
- projective dependency parsing (Eisner and Satta, 1999b)
- nonprojective dependency parsing (Cohen, Gómez-Rodríguez, et al., 2011)
- tree adjoining grammars (Alonso et al., 1999)
- combinatory categorical grammars (Kuhlmann and Satta, 2014)
- linear context-free rewriting systems (Burden and Ljunglöf, 2005)
- machine translation models (Melamed et al., 2004; Lopez, 2009)

We also mention Sikkel (1997), Gómez-Rodríguez (2010), and Goodman (1998), which include extensive examples of deduction systems related to syntactic analyses in NLP.

Deduction System Notation vs. Dyna Notation. The difference between semiring-weighted deduction systems and Dyna is largely superficial.

²¹⁵We also mention Teitelbaum (1973) who showed that the CKY algorithm could be extended to work in any semiring. Goodman (1998) acknowledges this paper as laying the foundation for most of the work on semiring-weighted deduction.

Consider a Dyna rule of the form:

$$x \oplus = y_1 \circ \cdots \circ y_K$$

It is rendered in deduction systems notation as

$$\frac{\boldsymbol{y}_1 \quad \cdots \quad \boldsymbol{y}_K}{\boldsymbol{x}}$$

Deduction notation is often used in research papers and books. However, it is not a user-friendly programming language.²¹⁶ Additionally, the Dyna notation is more explicit about how item values are combined.²¹⁷

Deductive Databases and Logic Programming. The Dyna notation is a natural generalization of Prolog (and Datalog), which allows items to have non-boolean values. Green, Karvounarakis, et al. (2007)—independently from Eisner, Goldlust, et al. (2005)—discovered the beauty and utility of attaching semiring values to Datalog rules. Their work, however, focused on data provenance²¹⁸ rather than dynamic programming. We have built on the logic programming literature in many ways: our semantics, execution, and source-to-source transformations have roots in the logic programming literature. These connections are discussed throughout the dissertation.

Einsum Notation. Dyna programs are similar to einsum notation, which is widely used in scientific computing libraries toolkits, such as numpy (Harris

²¹⁶Gómez-Rodríguez et al. (2009) built a system for interpreting unweighted rules written in notation that mimics the deduction systems notation but in plain text.

²¹⁷In Dyna <u>2</u> (discussed in §D), item values can be combined in arbitrary ways (i.e., by more than the two semiring operations), which would not be possible in the deduction notation.

²¹⁸Data provenance is a record trail that traces why certain database items exist. In effect, data provenance provides a derivation that an item belongs in the database.

et al., 2020), as a concise notation for tensor products.

Example 47. Below are example pairs of einsum code (in a % comment above) and its equivalent Dyna code.

```
607 % out<sub>0</sub> = einsum('i->', a)

608 out<sub>0</sub> += a(I).

609 % out<sub>1</sub> = einsum('i,i->i', a, b)

610 out<sub>1</sub>(I) += a(I) * b(I).

611 % out<sub>2</sub> = einsum('i,j->ij', a, b)

612 out<sub>2</sub>(I,J) += a(I) * b(J).

613 % out<sub>3</sub> = einsum('ii->', c)

614 out<sub>3</sub> += c(I,I).

615 % out<sub>4</sub> = einsum('ij,ji->ij', c, d)

616 out<sub>4</sub>(I,J) += c(I,J) * d(J,I).

617 % out<sub>5</sub> = einsum('ij,kl->ijkl', c, d)

618 out<sub>5</sub>(I,J,K,L) += c(I,J) * d(K,L).

619 % out<sub>6</sub> = einsum('ij,jk->ik', c, d)

620 out<sub>6</sub>(I,K) += c(I,J) * d(J,K).
```

Dyna notation extends einsum in many ways—the most important being that it supports interlocking and recursive equations. We hope that someday Dyna—or some variant thereof—will supplant einsum as the default domain-specific language for tensor manipulation.

Context-Free Grammars. It is interesting to relate Dyna programs to context-free grammars (CFGs). The mapping is as follows

- (i) terminal \Longrightarrow value
- (ii) nonterminal \Longrightarrow ground item
- (iii) $x \rightarrow y z \Longrightarrow x += y * z$
- (iv) parse tree \Longrightarrow derivation
- (v) the yield of a parse tree \Longrightarrow the value of a derivation
- (vi) the language generated (by a nonterminal) ⇒ the total value of all derivations (of that item)

The primary difference is that Dyna does not restrict the set of terminals and nonterminals to finite sets. Anecdotally, exploring the connection between Dyna programs and CFGs is how we discovered the left-corner transformation (§7.5.2) and the Newton solver (§5.2). Lastly, we note that our program transformations may be useful as grammar transformations, and our search method may be effective at reducing grammar size.

9.2 Future Work

Beyond semirings. In future work, exploring the extensions of the Dyna 1 language would be interesting.

- We discuss extension to Dyna 2 in §D.
- We also mention that a relaxation of semirings that allows nonassociative product operations (e.g., Gildea (2020)) might be interesting to explore. This would allow the item values to be trees rather than strings of values. With some attention to detail, we believe many of our results will likely extend to this more general framework, but we leave that exploration for future work.

Improving efficiency. Efficiency can be improved in many ways that we have not explored in this dissertation. We highlight a few below.

Code generation and compilation are crucial to efficiently running Dyna programs. Eisner, Goldlust, et al. (2005) provided a compiler for Dyna, but it is no longer actively developed. We hope to provide a similar compiler someday.

- There are optimizations for certain semirings that can lead to more efficient execution strategies. For example, superior semirings are amenable to A*like execution algorithms (Huang, 2008; Felzenszwalb and McAllester, 2007; Knuth, 1977). We want to explore these in future work.
- Prioritization function optimization for use in the prioritized fixpoint iteration algorithm (§5.1.3) could improve our execution times considerably as we will not require an SCC-decomposition pass. It may be possible to synthesize prioritization functions using SMT as in Bagnara et al. (2010) and Podelski and Rybalchenko (2004), which was applied to imperative programs. Synthesizing an effective prioritization heuristic could greatly improve our ability to synthesize efficient code.
- Data structure optimization can dramatically improve our ability to store and retrieve values at the innermost loops of our execution algorithms. The following data structure choices can have a dramatic effect on the execution time in practice: memory layouts (e.g., row-order or column-order layout of a dense array in memory), sparse vs. dense representations of relations (e.g., hash tables vs. arrays), and indexes on relations (including sorted order).²¹⁹
- It would also be interesting to explore the memory-efficient mixed-chaining algorithms of (Filardo and Eisner, 2012; Filardo, 2017).

²¹⁹Vieira, Francis-Landau, et al. (2017) describe a reinforcement-learning technique for dynamically favoring faster data structures during the course of execution. It uses a policy that dynamically adjusts how much data is stored in each of the competing data structures without duplicating data.

Improving transformation search. We provided a proof-of-concept search algorithm, which is not very sophisticated. The following are some ideas for how it might be improved:

- Decomposing the set of programs explored into subtasks that may depend
 on one another enables a kind of dynamic programming that should improve
 search efficiency because each subtask can be optimized recursively and combined into programs that efficiently compute the required output tasks.
- Vieira, Francis-Landau, et al. (2017) proposed dynamic optimization which
 adapts the solver to specific workloads as it executes using reinforcement learning rather than ahead of time as we do in §8. They also propose decomposition
 into subtasks to speed up learning.

Broadening the set of program transformations. There are still many reusable tricks employed by algorithms designers that have not yet been systematized into program transformations, and other works are exploring the efficient compilation of combinatorial algorithms. One notable recent work is PERPL (Chiang, McDonald, et al., 2023), a probabilistic programming language that can efficiently compile away certain kinds of "linear" data structures, such as stacks and lists. We want to port over their compilation technique as a source-to-source transformation for Dyna programs. It would allow us to recover useful speed-ups relevant to NLP algorithms (e.g., push-down automata and linear indexed grammars). In principle, the efficient compilation strategy used in PERPL can transform the following inefficient program for parsing with a

push-down automaton

```
621 pda([],0).
622 pda([X|Xs],K) += pda(Xs,J) * rewrite(X,Y) * word(Y,J,K).
623 pda([X|Xs],K) += pda([Z,Y|Xs],K) * rewrite(X,Y,Z).
624 goal += length(N) * pda([s],N).
625
626 input: rewrite(_,_,); rewrite(_,_); length(_); word(_,_,_).
627 output: goal.
```

into an efficient algorithm that is equivalent to CKY.

```
628 p(I,X,K) += rewrite(X,Y) * word(Y,I,K).
629 p(I,X,K) += pda(I,Y,J) * p(J,Z,K) * rewrite(X,Y,Z).
630 goal += length(N) * p(0,s,N).
```

We discuss this example further in §E.12.

Chapter 10

Conclusion

We presented Dyna as a notation for dynamic programming (§2). We showed in Example 1 and Example 2 that NLP authors have traditionally overlooked speedups. We argued that this demonstrates a need for tools to facilitate the analysis and development of NLP algorithms.

We studied Dyna's foundations: semantics (§ 4) and execution (§ 5). We provided a new sum-over-derivations semantics. We discovered a new execution algorithm, Newton, and we generalized Eisner, Goldlust, et al. (2005)'s forward-chaining algorithm to support non-range-restricted programs. We identified and explored several rich families of source-to-source transformations that systematize the speedup tricks found across many papers for deriving more efficient algorithms, including the generalized fold/unfold, speculation, and left-corner transformations.

We provided a search-based system (§8) for finding a sequence of meaningpreserving source-to-source transformations to a given initial program that minimizes its worst-case running time. Our program analyzer (§6) can automate much of the algorithmic analysis found in NLP papers, including deriving precise big- \mathcal{O} bounds on runtime and space complexity. We showed that our system recovers many interesting oversights published in the NLP literature, including oversights in the analysis (Example 2) and overlooked speedups (Example 1).

We have released a freely available prototype implementation https://github.com/timvieira/dyna-pi. We have found that the utility of our system is that it aids in discovering the right (i.e., efficient) high-level algorithm. In other words, our tool is most helpful to algorithms developers and the authors of NLP systems in finding the right recurrence—similar to a theorem-proving assistant. Our implementation does not yet perform efficient code generation as Eisner, Goldlust, et al. (2005). The implementation includes several useful development utilities, such as visualization of derivations, that we have not discussed here.

We believe that this dissertation represents a significant step toward automating the development of efficient dynamic programs, especially those found in natural language processing applications.

Appendix A

Inference Algorithms

This section shows how our algorithms can solve certain programs with infinitely many non-0 items.

Example 48. Below, we have three versions of the geometric series program. On the left, we have the basic geometric series. In the middle, we have threaded through a useless variable 1. On the right, we add an inequality constraint on 1. Beneath each program variant, we show their respective charts on iterations 1, 2, 3, and the final iteration.

Example 49 (Constraint Accumulation). Our method does not, however, solve

$$\mathfrak{p} = \begin{cases} x(I) += (I > 0). \\ x(I) += 0.5 * (I > J) * x(J). \end{cases}$$

After five iterations, we see that the constraints are just accumulating,

```
\mathfrak{m}^{(5)} = \begin{cases} x(I) \mathrel{+=} (I > \emptyset). \\ x(I) \mathrel{+=} \emptyset.5 * (I > J) * (J > \emptyset). \\ x(I) \mathrel{+=} \emptyset.25 * (I > J) * (J > J_1) * (J_1 > \emptyset). \\ x(I) \mathrel{+=} \emptyset.125 * (I > J) * (J > J_1) * (J_1 > J_2) * (J_2 > \emptyset). \\ x(I) \mathrel{+=} \emptyset.0625 * (I > J) * (J > J_1) * (J_1 > J_2) * (J_2 > J_3) * (J_3 > \emptyset). \end{cases}
```

To solve this program, we need to simplify the rule bodies. Specifically, we need a method that counts the number of satisfying assignments to the delayed constraints and projects any constraints on the variables in the head. For example, to simplify the bottom rule, we need an efficient mechanism for representing the solution to

```
multiplicity(I) += (I > J) * (J > J_1) * (J_1 > J_2) * (J_2 > J_3) * (J_3 > \emptyset)
```

Efficiently representing and reasoning about the solutions to the multiplicity problems encountered during fixpoint iteration in a general way is challenging. It is possible that SMT solvers (e.g., Moura and Bjørner (2008)) may be helpful for certain constraint domains like arithmetic inequalities. We think this is an interesting area for future work.

Example 50. Consider the program for computing the reflexive, transitive closure of a graph in the min-plus semiring:

```
path(I,I) min= 0
path(I,K) min= path(I,J) + edge(J,K).
```

Notice that the top rule is not range-restricted, i.e., it has a free variable in its head. As discussed in §5.1.3, this kind of free variable will cause problems for

Eisner, Goldlust, et al. (2005)'s forward chaining algorithm.

Below, we show naïve fixpoint iteration on the following input data:

```
edge(a,b) min= 1.
edge(b,c) min= 1.
edge(c,a) min= 1.
```

The chart of each iteration is shown below:

```
\mathfrak{m}^{(1)} = \begin{cases} \text{edge}(\mathtt{a},\mathtt{b}) & \text{min= 1.} \\ \text{edge}(\mathtt{b},\mathtt{c}) & \text{min= 1.} \\ \text{edge}(\mathtt{c},\mathtt{a}) & \text{min= 1.} \\ \text{edge}(\mathtt{c},\mathtt{a}) & \text{min= 1.} \\ \text{path}(\mathtt{I},\mathtt{I}) & \text{min= 0.} \end{cases} \\ \mathfrak{m}^{(2)} = \begin{cases} \text{edge}(\mathtt{a},\mathtt{b}) & \text{min= 1.} \\ \text{edge}(\mathtt{c},\mathtt{a}) & \text{min= 1.} \\ \text{path}(\mathtt{I},\mathtt{I}) & \text{min= 0.} \\ \text{path}(\mathtt{a},\mathtt{b}) & \text{min= 1.} \\ \text{path}(\mathtt{a},\mathtt{b}) & \text{min= 1.} \\ \text{path}(\mathtt{a},\mathtt{b}) & \text{min= 1.} \\ \text{path}(\mathtt{a},\mathtt{b}) & \text{min= 2.} \\ \text{path}(\mathtt{b},\mathtt{c}) & \text{min= 2.} \\ \text{path}(\mathtt{b},\mathtt{c}) & \text{min= 2.} \\ \text{path}(\mathtt{b},\mathtt{c}) & \text{min= 2.} \end{cases} \\ \mathfrak{m}^{(2)} = \begin{cases} \text{edge}(\mathtt{a},\mathtt{b}) & \text{min= 1.} \\ \text{edge}(\mathtt{c},\mathtt{a}) & \text{min= 1.} \\ \text{path}(\mathtt{b},\mathtt{c}) & \text{min= 1.} \\ \text{path}(\mathtt{c},\mathtt{a}) & \text{min= 2.} \\ \text{path}(\mathtt{b},\mathtt{a}) & \text{min= 2.} \\ \text{path}(\mathtt{c},\mathtt{a}) & \text{min= 3.} \\ \text{path}(\mathtt{b},\mathtt{b}) & \text{min= 3.} \\ \text{path}(\mathtt{b},\mathtt{b}) & \text{min= 3.} \\ \text{path}(\mathtt{b},\mathtt{b}) & \text{min= 3.} \\ \text{path}(\mathtt{c},\mathtt{c}) & \text{
```

We see that fixpoint iteration converges after 4 iterations. The thing to notice about these charts is that they contain a nonground rule path(I,I) min= 0. What that means is when we query for the value of path(a,a), for example, two rules in $m^{(4)}$ contribute to its value giving min(0,3) = 0.

Appendix B

Program Transformations

B.1 Strong Equivalence

Proposition 5. If $\mathfrak{p}^{(s)} \stackrel{\text{strong}}{\equiv} \mathfrak{p}^{(t)}$, then $\mathfrak{p}^{(s)} \equiv \mathfrak{p}^{(t)}$.

Proof. Fix an arbitrary item $x \in \mathbb{H}$.

$$\mathfrak{p}^{(s)}\llbracket x \rrbracket = \bigoplus_{\boldsymbol{\delta}^{(s)} \in \mathfrak{p}^{(s)} \{x\}} \mathfrak{p}^{(s)} \llbracket \boldsymbol{\delta}^{(s)} \rrbracket \qquad \qquad \text{[Definition 22]} \qquad \text{(B.1)}$$

$$= \bigoplus_{\boldsymbol{\delta}^{(s)} \in \mathfrak{p}^{(s)} \{x\}} \mathfrak{p}^{(t)} \llbracket \boldsymbol{\phi}(\boldsymbol{\delta}^{(s)}) \rrbracket \qquad \qquad \text{[value preserving]} \qquad \text{(B.2)}$$

$$= \bigoplus_{\boldsymbol{\delta}^{(t)} \in \mathfrak{p}^{(t)} \{x\}} \mathfrak{p}^{(t)} \llbracket \boldsymbol{\delta}^{(t)} \rrbracket \qquad \qquad \text{[bijection]} \qquad \text{(B.3)}$$

$$= \mathfrak{p}^{(t)} \llbracket x \rrbracket \qquad \qquad \text{[Definition 22]} \qquad \text{(B.4)}$$

B.2 Correctness of the Fold and Unfold

Transformations

This appendix provides several results that characterize the correctness of the {basic, generalized} \times {unfold, fold} transformations. We have chosen to give the most detailed correctness proofs for the case of generalized unfold since the proof of correctness for the generalized fold transformation follows the same argument. Furthermore, the correctness proofs for basic unfold and fold transformations (Theorem 4) and Theorem 5) are corollaries of the generalized unfold proof. The correctness proof for the generalized unfold transformation is given by two lemmas, Lemma 1 and Lemma 4, that are assembled into the main theorem: Theorem 2. Lemma 1 provides generalize conditions under which the derivation mapping ϕ (provided in Algorithm 24) is a value-preserving mapping. Then, Lemma 4 demonstrates that ϕ (provided in Algorithm 32) is the inverse of ϕ . Analogous results are provided for the generalized fold transformation. We remind the reader that the conditions provided here are not easy to verify (§7.4.2.3). A set of practical sufficient conditions are provided in §7.4.2.4.

Overview. An overview of the results is given below.

	unfold	fold
Derivation mapping (ϕ) $_{s o t}$		
\subseteq Algorithm	Algorithm 24	Algorithm 25
$\hookrightarrow \mathop{\phi}\limits_{s o t}$ is a value-preserving mapping	Lemma 1	Lemma 5
Inverse derivation mapping (ϕ) $_{t o s}$		
\subseteq Algorithm	Algorithm 31	Algorithm 32
$\hookrightarrow \underset{t \to s}{\phi}$ is the inverse of $\underset{s \to t}{\phi}$	Lemma 4	Lemma 6
Semantics-preserving		
\hookrightarrow basic	Theorem 4	Theorem 5
\subseteq generalized (impractical)	Theorem 2	Theorem 3
\subseteq generalized (practical)	Theorem 1	

Additional definitions. Some of our proofs make use of the following ordering relation on derivations:

Definition 43. We define the **subtree-ordering relation**. Let $\delta, \delta' \in \mathfrak{p}\{\cdot\}$ be derivations (of possibly different items). We say that $\delta' \prec \sup_{subtree} \delta$ if and only if δ' is a (strict) subtree of δ .

Lemma 1. Suppose

$$\begin{bmatrix} \mathsf{A}_2 \end{bmatrix} \ \mathfrak{p}^{(s)} \overset{\text{strong}}{\equiv} \mathfrak{p}^{(d)}$$

A₃ the method $\phi_{s \to t}: \mathfrak{p}^{(s)}\{\cdot\} \to \mathfrak{p}^{(t)}\{\cdot\}$ specified in Algorithm 24 always terminates.

Then, $\phi \atop s \to t$ satisfies the following, $\forall \pmb{\delta}^{(s)} \in \mathfrak{p}^{(s)}\{\cdot\}$:

where x is the root label of $\delta^{(s)}$.

Proof. Our correctness argument uses induction on the ordering $\underset{\text{steps}}{\prec}$ over $\mathfrak{p}^{(s)}\{\cdot\}$.

Since ϕ terminates by A_3 and it is a side-effect-free deterministic procedure, there exists a finite number of steps per call to ϕ . Let $\operatorname{steps}\left\{\phi\left(\delta\right)\right\}$ be the number of steps when ϕ is run on δ . The ordering relation \prec is defined as follows: $\delta' \prec \delta \iff \operatorname{steps}\left\{\phi\left(\delta'\right)\right\} < \operatorname{steps}\left\{\phi\left(\delta'\right)\right\}$. We may assume that every function call has a positive overhead cost so that no subroutine call can have zero cost; this ensures that \prec is irreflexive. Clearly, if $\phi\left(\delta\right)$ recursively calls $\phi\left(\delta'\right)$ (i.e., $\phi\left(\delta'\right)$ is a subroutine of $\phi\left(\delta\right)$), and $\phi\left(\delta\right)$ holds, then $\delta' \prec \delta$. Furthermore, $\phi\left(\delta'\right)$ is a well-founded ordering on $\phi\left(\delta\right)$ steps

 $\phi_{s \to t}(\cdot)$ because there are no infinite descending chains of comparison. Thus, \prec is safe for use in an inductive proof using the principle of Noetherian induction (Cohn, 1969).

Choose $\delta^{(s)} \in \mathfrak{p}^{(s)}\{\cdot\}$ arbitrarily.

Inductive hypothesis:

$$\boxed{\mathsf{IH}} \ \forall \boldsymbol{\delta}' \underset{\mathsf{steps}}{\prec} \boldsymbol{\delta}^{(s)} : \underset{s \to t}{\phi}(\boldsymbol{\delta}') \ \mathsf{exists}, \ \underset{s \to t}{\phi}(\boldsymbol{\delta}') \in \mathfrak{p}^{(t)}\{\boldsymbol{x}\}, \ \mathsf{and} \ \llbracket \boldsymbol{\delta}' \rrbracket = \left[\begin{smallmatrix} \phi \\ s \to t \end{smallmatrix} \right].$$

Base case: ($\delta^{(s)}$ is a constant or builtin): This case is straightforward.

 C_2 : $\phi_{s \to t}(\delta^{(s)}) \in \mathfrak{p}^{(t)}\{\cdot\}$ because constants and builtins are, by Definition 23, always included in the set of derivations.

Inductive case: There are three cases to consider:

Case 1 ($\delta^{(s)}$ is a product $\overset{y_1}{\triangle} \circ \cdots \circ \overset{y_K}{\triangle}$): This case is straightforward. Notice $\overset{y_1}{\triangle}, \ldots, \overset{y_K}{\triangle} \prec \delta^{(s)}$ because they are subroutines.

$$\begin{array}{c} \boxed{\textbf{C}_1} \text{:} \quad \phi \left(\overset{y_1}{\bigtriangleup} \circ \cdots \circ \overset{y_K}{\bigtriangleup} \right) \stackrel{\text{def}}{=} \quad \phi \left(\overset{y_1}{\bigtriangleup} \right) \circ \cdots \circ \phi \left(\overset{y_K}{\bigtriangleup} \right) \end{aligned}$$

C₂:

$$\begin{array}{c} \overset{y_1 \circ \cdots \circ y_K}{\bigtriangleup} \in \mathfrak{p}^{(s)} \{ y_1 \circ \cdots \circ y_K \} \Longrightarrow \overset{y_1}{\bigtriangleup} \in \mathfrak{p}^{(s)} \{ y_1 \}, \ldots, \overset{y_K}{\bigtriangleup} \in \mathfrak{p}^{(s)} \{ y_K \} \\ & \stackrel{\text{IIH}}{\Longrightarrow} \phi \left(\overset{y_1}{\bigtriangleup} \right) \in \mathfrak{p}^{(t)} \{ y_1 \}, \ldots, \phi \left(\overset{y_K}{\bigtriangleup} \right) \in \mathfrak{p}^{(t)} \{ y_K \} \\ & \Longrightarrow \phi \left(\overset{y_1}{\bigtriangleup} \circ \cdots \circ \overset{y_K}{\bigtriangleup} \right) \in \mathfrak{p}^{(t)} \{ y_1 \circ \cdots \circ y_K \} \end{array}$$

$$\boxed{\textbf{C}_{3}} \text{:} \quad \left[\!\!\left[\begin{matrix} \phi \\ \\ \\ s \rightarrow t \end{matrix} \begin{pmatrix} \textbf{y}_{1} \\ \\ \\ \end{matrix} \circ \cdots \circ \begin{matrix} \textbf{y}_{K} \\ \\ \end{matrix} \end{pmatrix} \right] \right] = \left[\!\!\left[\begin{matrix} \phi \\ \\ \\ \\ \end{matrix} \end{pmatrix} \begin{pmatrix} \textbf{y}_{1} \\ \\ \\ \end{matrix} \end{pmatrix} \circ \cdots \circ \begin{matrix} \phi \\ \\ \\ \end{matrix} \end{pmatrix} \right] \left[\!\!\left[\begin{matrix} \textbf{H} \\ \\ \\ \end{matrix} \right] \circ \cdots \circ \begin{matrix} \textbf{y}_{K} \\ \\ \\ \end{matrix} \right] \right]$$

Case 2 (top rule of $\delta^{(s)}$ is in $\mathfrak{p}^{(t)}$): This case is straightforward. Observe that $\bigwedge_{\text{steps}}^{\mu} \delta^{(s)}$ because it is a subroutine.

$$\begin{array}{c|c}
\hline
C_1: & \phi \\ \downarrow \\ s \to t \\
\hline
\begin{pmatrix} x \\ \downarrow \\ \mu \\ \\ \end{matrix}
\end{pmatrix} \stackrel{\text{def}}{=} \begin{pmatrix} x \\ \downarrow \\ \\ \phi \\ \\ x \to t \\
\end{matrix}$$

 $\boxed{ \textbf{C}_2 } \text{: Since, the top rule of } \boldsymbol{\delta}^{(t)} = \underset{s \to t}{\phi} \left(\boldsymbol{\delta}^{(s)} \right) \text{ is a specialization of a rule in } \mathfrak{p}^{(t)} \\ \text{and its subtrees are in } \mathfrak{p}^{(t)} \{ \cdot \} \text{ by } \boxed{\textbf{H}} \text{, then } \boldsymbol{\delta}^{(t)} \in \mathfrak{p}^{(t)} \{ \boldsymbol{x} \}.$

Case 3 (top rule of $\delta^{(s)}$ is not in $\mathfrak{p}^{(t)}$): By Definition 35, the derivation $\delta^{(s)}$ must have the form:

$$oldsymbol{\delta}^{(s)} = \overbrace{igwedge}^{oldsymbol{x}} oldsymbol{\mathcal{Y}}_{j} \quad oldsymbol{eta}$$

We now unfold the derivation $\overset{y_j}{\triangle}$: $\left\langle \overset{\mu}{\triangle}, \mathbf{T}^{(\ell)} \right\rangle = u \begin{pmatrix} \overset{y_j}{\triangle} \end{pmatrix}$. This step is required for the (distinct) replacement rule $\mathbf{T}^{(\ell)} \in \mathcal{S}$ to fire in the transformed derivation (below). Observe that $\overset{\alpha}{\triangle}, \overset{\mu}{\triangle}, \overset{\beta}{\triangle} \preccurlyeq \delta^{(s)}$ because they are subroutines.

 C_3 :

$$\begin{bmatrix} \phi & x \\ & & & \\ &$$

Lemma 2. Suppose u is called by Algorithm 24 with $\mathfrak{p}^{(s)} \xrightarrow{\mathsf{unfold} \frac{d,j}{r \to \mathcal{S}}} \mathfrak{p}^{(t)}$ and

$$oldsymbol{\delta}^{(s)} = oldsymbol{\gamma}_{y_j} oldsymbol{eta}_{igwedge}$$

where the top rule of $\delta^{(s)}$ is the rule to be unfolded r. Then, $\left\langle \stackrel{\mu}{\bigtriangleup}, \mathbf{T}^{(\ell)} \right\rangle \leftarrow u \left(\stackrel{y_j}{\bigtriangleup} \right)$ has the following properties:

$$\boxed{\mathsf{C}_1} \ \bigwedge^{\mu} \in \mathfrak{p}^{(s)}\{\cdot\}$$

$$\begin{bmatrix}
C_2
\end{bmatrix}^{\mu} = \begin{bmatrix}
y_j\\
\triangle
\end{bmatrix}$$

 $oxed{\mathsf{C}_3} \ \mathbf{T}^{(\ell)} \in \mathfrak{p}^{(t)} \ and \ \mathbf{T}^{(\ell)} \ has \ the \ form \ oldsymbol{x} \oplus oldsymbol{=} oldsymbol{lpha} \circ oldsymbol{eta}$

Proof sketch. Conclusions C_1 and C_2 are obvious consequences of $\mathfrak{p}^{(s)} \stackrel{\text{strong}}{\equiv} \mathfrak{p}^{(d)}$ (assumed by the calling context).

C₃: By the construction of the unfold transformation (Definition 35), we know that r has the form $x \oplus = \alpha \circ y_j \circ \beta$. By inspecting Definition 35, we can see that the rule $\mathbf{T}^{(\ell)}$ must be of the form $x \oplus = \alpha \circ \mu \circ \beta$ for some μ . The transformation replaces the item y_j with the various bodies of the defining rules from $\mathfrak{p}^{(d)}$. Furthermore, when we know the definition rule $\mathbf{D}^{(\ell)}$ that was used, we know which rule in $\mathfrak{p}^{(t)}$ replaced r. Specifically, $\mathbf{T}^{(\ell)} = \frac{\sigma}{d \to t}(\mathbf{D}^{(\ell)})$ using the one-to-one alignment $\frac{\sigma}{d \to t}$ (described in Definition 35). Any variable bindings have already been propagated when unfold created $\mathbf{T}^{(\ell)}$.

Lemma 3. The method u^{-1} correctly implements the inverse of u.

Proof (Sketch). Our argument is that u^{-1} inverts u by working backwards from the output of \mathbf{u} and inverting each line one-at-a-time. 220

Suppose the calling context has $\boldsymbol{\delta}^{(s)}$ with top rule $\mathbf{T}^{(\ell)} \in \mathcal{S}$. Thus, $\mathbf{T}^{(\ell)}$ has the form $x \oplus = \alpha \circ \mu \circ \beta$ and the unfolded rule r has the form $x \oplus = \alpha \circ y_j \circ \beta$.

The only tricky part of inverting u is that we have to reconstruct the intermediate item y_j that u erased on line 24.13. Fortunately, there is just enough information in the transformation to reconstruct it from the identity of $\mathbf{T}^{(\ell)}$. The generalized unfold transformation (Definition 35) is careful to ensure this is possible by tracking the alignment $\sigma_{x_{d-1}}$.

The way that line 25.15 manages to reconstruct the erased item is as follows: We look up the distinct unfolder rule $\mathbf{D}^{(\ell)} \in \mathfrak{p}^{(d)}$ that was used to create $\mathbf{T}^{(\ell)}$ during the generalized unfold transformation (Definition 35). We have $\mathbf{D}^{(\ell)} = \frac{\sigma}{t \to d}(\mathbf{T}^{(\ell)})$. Furthermore, $\mathbf{T}^{(\ell)}$ must be of the form $y' \oplus = \mu'$. Now, an intermediate subtree must exist (that unfold erased on line 24.13), built from an instantiation of the rule $\mathbf{D}^{(\ell)}$. To build this node, we determine the substitution that makes the folder's body $\mu' \theta = \mu$. We can determine a unique substitution $\theta = \operatorname{cover}(\mu', \mu)$. We can now determine the name y_j of the missing intermediate item $y_j = \theta(y')$.

The remaining operations are correct by the assumption that $\mathfrak{p}^{(s)} \stackrel{\text{strong}}{=} \mathfrak{p}^{(d)}$.

 $[\]overline{^{220}\text{Recall that if }h(x)=f(g(x)), \text{ then }h^{-1}(x)=g^{-1}(f^{-1}(x)) \text{ (assuming these inverses exist).}$

Algorithm 31 Inverse derivation mapping for $\mathfrak{p}^{(s)} \xrightarrow{\text{unfold } r \to \mathcal{S}} \mathfrak{p}^{(t)}$.

- 1. \triangleright The comments at the top of Algorithm 24 apply here as well.
- 2. $\operatorname{def} \phi_{t \to s} \left(\delta^{(t)} \right)$:

 \triangleright Inverse derivation mapping

- 3. **case** constant or builtin: **return** $\boldsymbol{\delta}^{(t)}$
- 4. **case** $\bigwedge^{y_1} \circ \cdots \circ \bigwedge^{y_K}$: **return** $\phi \left(\bigwedge^{y_1} \right) \circ \cdots \circ \phi \left(\bigwedge^{y_K} \right)$
- 5. **case** $\begin{pmatrix} x \\ | \\ \wedge \\ \end{pmatrix}$ **and** its top rule $\in \mathfrak{p}^{(s)}$: **return** $\phi \atop t \to s} \begin{pmatrix} \mu \\ \wedge \\ \end{pmatrix}$
- 6. **case** $\begin{pmatrix} x \\ \alpha & \mu & \beta \\ \triangle & \triangle \end{pmatrix}$ **and** its top rule $\notin \mathfrak{p}^{(s)}$:

riangleright Top rule must be in S

- 7. **return** $\phi \left(\bigwedge^{\alpha} \right) \quad u^{-1} \left(\phi \left(\bigwedge^{\mu} \right) \right) \quad \phi \left(\bigwedge^{\beta} \right)$ $t \to s \quad t \to t \to s \quad t \to s \quad t \to$
- $\begin{tabular}{l} \rhd \textit{The top rule of the transformed derivation is } r \end{tabular}$

Algorithm 32 Inverse derivation mapping for $\mathfrak{p}^{(s)} \xrightarrow{\text{fold } \mathcal{S} \to r} \mathfrak{p}^{(t)}$.

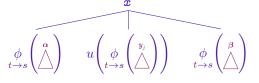
- 1. \triangleright The comments at the top of Algorithm 24 apply here as well.
- 2. $\operatorname{\mathbf{def}} \phi_{t \to s} \left(\boldsymbol{\delta}^{(t)} \right)$:

▷ Inverse derivation mapping

- case constant or builtin: return $\boldsymbol{\delta}^{(t)}$
- 4. **case** $\bigwedge^{y_1} \circ \cdots \circ \bigwedge^{y_K}$: **return** $\phi \left(\bigwedge^{y_1} \right) \circ \cdots \circ \phi \left(\bigwedge^{y_K} \right)$
- 5. **case** $\begin{pmatrix} x \\ | \\ \wedge \\ \end{pmatrix}$ **and** its top rule $\in \mathfrak{p}^{(s)}$: **return** $\phi \in \mathcal{P}(s)$
- 6. **case** $\begin{pmatrix} x \\ x \\ y_j & \beta \end{pmatrix}$ **and** its top rule $\notin \mathfrak{p}^{(s)}$:

 \triangleright Top rule must be r

7. return



ightharpoonup The top rule of the transformed derivation is given by the call to u

Lemma 4. Suppose

$$\boxed{ \mathsf{A}_1 } \ \mathfrak{p}^{(s)} \xrightarrow{\mathsf{unfold} \frac{d,j}{r \to \mathcal{S}}} \mathfrak{p}^{(t)}$$

$$\boxed{\mathsf{A}_2} \ \mathfrak{p}^{(s)} \overset{\mathrm{strong}}{\equiv} \mathfrak{p}^{(d)}$$

 $igotimes_{s o t}: \mathfrak{p}^{(s)}\{\cdot\} o \mathfrak{p}^{(t)}\{\cdot\} \ is \ the \ partial^{221} \ function \ specified \ in \ Algorithm \ 24.$

Then, $\forall \boldsymbol{\delta}^{(t)} \in \mathfrak{p}^{(t)}\{\cdot\} : \exists ! \, \boldsymbol{\delta}^{(s)} \in \mathfrak{p}^{(s)}\{\cdot\} : \oint\limits_{s \to t} \left(\boldsymbol{\delta}^{(s)}\right) = \boldsymbol{\delta}^{(t)}$. Furthermore, the function $\phi_{t \to s}$ specified in Algorithm 31 computes $\boldsymbol{\delta}^{(s)} = \phi_{t \to s}\left(\boldsymbol{\delta}^{(t)}\right)$.

Proof. Our proof uses induction on \prec over $\mathfrak{p}^{(t)}\{\cdot\}$. 222 Fix an arbitrary $\boldsymbol{\delta}^{(t)} \in \mathfrak{p}^{(t)}\{\cdot\}$.

Inductive hypothesis:

Base case: ($\delta^{(t)}$ is a constant or builtin): This case is straightforward.

We seek to identify all $\delta^{(s)}$ that satisfy the following equation:

$$\boldsymbol{\delta}^{(t)} = \phi \left(\boldsymbol{\delta}^{(s)} \right)$$

It is clear from the implementation of $\phi_{s \to t}$ (see line 24.5) that $\pmb{\delta}^{(t)}$ has this form is if and only if

$$\boldsymbol{\delta}^{(t)} = \boldsymbol{\delta}^{(s)}$$

²²¹Here, we take nontermination on a given input to mean that the function is undefined on it. ²²²The fact that we call pull off the induction on an ordering (\prec) that is guaranteed to be well-founded without additional assumptions is why ϕ is a total function and why ϕ is only guaranteed to be a partial function. The induction scheme used in the proof of Lemma 1 required us to stipulate that ϕ terminates to perform induction in the number of steps performed by Algorithm 24.

This is clearly a unique solution. Furthermore, $\boldsymbol{\delta}^{(s)} = \phi_{t \to s} \left(\boldsymbol{\delta}^{(t)} \right) = \boldsymbol{\delta}^{(t)}$.

<u>Inductive case</u>: As one might expect, the construction of $\phi_{t \to s}$ essentially inverts each step of the construction of $\phi_{s \to t}$. Thus, there are three analogous cases to consider as in Lemma 1:

Case 1 ($\delta^{(t)}$ is a product $\gamma_1 \circ \cdots \circ \gamma_K$): This case is straightforward.

We seek to identify all $\delta^{(s)}$ that satisfy the following equation:

$$oldsymbol{\delta}^{(t)} = igwedge_{1}^{y_{1}} \circ \cdots \circ igwedge_{N}^{y_{K}} = igoplus_{s o t} igg(oldsymbol{\delta}^{(s)}igg)$$

It is clear from the implementation of $\phi_{s \to t}$ (see line 24.6) that $\pmb{\delta}^{(t)}$ has this form is if and only if

Now, by $\overline{\mathbb{IH}}$, we have that for each $k \in [1:K]: \exists ! \, \boldsymbol{\delta}_k^{(s)}: \overset{y_k}{\triangle} = \underset{s \to t}{\phi} \left(\boldsymbol{\delta}_k^{(s)}\right);$ and, furthermore, that $\boldsymbol{\delta}_k^{(s)} = \underset{t \to s}{\phi} \left(\overset{y_k}{\triangle}\right).$ Therefore, the equation $\boldsymbol{\delta}^{(t)} = \underset{s \to t}{\phi} \left(\boldsymbol{\delta}^{(s)}\right)$ has exactly one solution for $\boldsymbol{\delta}^{(s)}$. Furthermore,

$$\boldsymbol{\delta}^{(s)} = \boldsymbol{\delta}_1^{(s)} \circ \cdots \circ \boldsymbol{\delta}_K^{(s)} = \phi \left(\bigwedge_{t \to s}^{y_1} \right) \circ \cdots \circ \phi \left(\bigwedge_{t \to s}^{y_K} \right) = \phi \left(\boldsymbol{\delta}^{(t)} \right)$$

Case 2 (top rule of $\delta^{(t)} \in \mathfrak{p}^{(s)}$): This case is straightforward.

We seek to identify all $\delta^{(s)}$ that satisfy the following equation:

$$oldsymbol{\delta}^{(t)} = egin{array}{c} oldsymbol{x} \\ ig| \\ igwedge \\ & igwedge \\ \end{array} = oldsymbol{\phi}_{s o t} igg(oldsymbol{\delta}^{(s)}igg)$$

It is clear from the implementation of $\phi_{s\to t}$ (see line 24.7) that $\pmb{\delta}^{(t)}$ has this form is if and only if

Now, by $\overline{\mathbb{IH}}$, we have that $\exists ! \delta_1^{(s)} : \overset{\mu}{\triangle} = \underset{s \to t}{\phi} \left(\delta_1^{(s)} \right)$; and, furthermore, that $\delta_1^{(s)} = \underset{t \to s}{\phi} \left(\overset{\mu}{\triangle} \right)$. From line 24.7, it is clear that the top rule must be the same in $\delta^{(s)}$ and $\delta^{(t)}$; therefore, the equation $\delta^{(t)} = \underset{s \to t}{\phi} \left(\delta^{(s)} \right)$ has exactly one solution for $\delta^{(s)}$. Furthermore,

Case 3 (top rule of $\delta^{(t)} \notin \mathfrak{p}^{(s)}$):

We seek to identify all $\delta^{(s)}$ that satisfy the following equation:

$$\boldsymbol{\delta}^{(t)} = \overbrace{\alpha \quad \mu \quad \beta}^{x} = \phi \left(\boldsymbol{\delta}^{(s)}\right)$$

It is clear from the implementation of $\phi_{s \to t}$ (see line 24.9) that $\pmb{\delta}^{(t)}$ has this form is if and only if

$$\begin{array}{cccc}
x & x \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & &$$

By IH, we have unique solutions to the following equations:

Furthermore, those solutions are given by ϕ :

$$\phi_{t \to s} \begin{pmatrix} \alpha \\ \triangle \end{pmatrix} = \delta_0^{(s)} \quad ; \quad \phi_{t \to s} \begin{pmatrix} \mu \\ \triangle \end{pmatrix} = \delta_1^{(s)} \quad ; \quad \phi_{t \to s} \begin{pmatrix} \beta \\ \triangle \end{pmatrix} = \delta_2^{(s)}$$

Next, we solve $\delta_1^{(s)} = u\left(\delta_3^{(s)}\right)$. According to Lemma 3, this equation has a unique solution that is computed by $\delta_3^{(s)} = u^{-1}\left(\delta_1^{(s)}\right)$ with line 25.7. Lastly, we have that the unique solution $\delta^{(s)}$ is given by

$$\boldsymbol{\delta}^{(s)} = \phi \begin{pmatrix} x \\ \alpha & \mu & \beta \\ \triangle & \triangle & \triangle \end{pmatrix} \stackrel{\text{def}}{=} \psi \begin{pmatrix} \alpha \\ \alpha \\ \Delta \end{pmatrix} u^{-1} \begin{pmatrix} \phi \\ \alpha \\ \Delta \end{pmatrix} \psi \begin{pmatrix} \beta \\ \Delta \end{pmatrix}$$

255

Theorem 2. Given $\mathfrak{p}^{(s)} \xrightarrow{\text{unfold} \frac{d,j}{r \to S}} \mathfrak{p}^{(t)}$ and $\mathfrak{p}^{(s)} \stackrel{\text{strong}}{=} \mathfrak{p}^{(d)}$. Suppose the function $\phi \left(\delta^{(s)} \right)$, as specified in Algorithm 24, terminates for all $\delta^{(s)} \in \mathfrak{p}^{(s)} \{\cdot\}$. Then, $\phi \in \mathfrak{p}^{(s)} \{\cdot\}$ is a value-preserving bijection between $\mathfrak{p}^{(s)} \{\cdot\}$ and $\mathfrak{p}^{(t)} \{\cdot\}$. Furthermore, $\mathfrak{p}^{(s)} \stackrel{\text{strong}}{=} \mathfrak{p}^{(t)}$.

Proof. Under the above conditions, it follows from Lemma 1 and Lemma 4 that ϕ is a value-preserving bijection between $\mathfrak{p}^{(s)}\{\cdot\}$ and $\mathfrak{p}^{(t)}\{\cdot\}$. Such a bijection ensures that $\mathfrak{p}^{(s)} \stackrel{\text{strong}}{\equiv} \mathfrak{p}^{(t)}$ by Definition 31.

Lemma 5. Suppose

$$\boxed{\mathsf{A}_2} \ \mathfrak{p}^{(s)} \stackrel{\mathrm{strong}}{\equiv} \mathfrak{p}^{(d)}$$

A₃ the method $\phi_{s \to t}: \mathfrak{p}^{(s)}\{\cdot\} \to \mathfrak{p}^{(t)}\{\cdot\}$ specified in Algorithm 25 always terminates.

Then, $\phi_{s \to t}$ satisfies the following, $\forall \delta^{(s)} \in \mathfrak{p}^{(s)}\{\cdot\}$:

$$oxed{\mathsf{C_1}} \oint\limits_{s o t} \left(oldsymbol{\delta}^{(s)}
ight) \textit{exists}$$

$$\begin{bmatrix}
\mathsf{C}_3 & \left[\phi \left(\boldsymbol{\delta}^{(s)} \right) \right] = \left[\boldsymbol{\delta}^{(s)} \right].$$

where x is the root label of $\delta^{(s)}$.

Proof. Our correctness argument uses induction on the ordering \prec over $\mathfrak{p}^{(s)}\{\cdot\}$ that we defined in the proof of Lemma 1.

Choose $\delta^{(s)} \in \mathfrak{p}^{(s)}\{\cdot\}$ arbitrarily.

Inductive hypothesis:

$$\boxed{\mathsf{IH}} \ \forall \pmb{\delta}' \underset{\mathsf{steps}}{\prec} \pmb{\delta}^{(s)} : \underset{s \to t}{\phi}(\pmb{\delta}') \ \mathsf{exists}, \ \underset{s \to t}{\phi}(\pmb{\delta}') \in \mathfrak{p}^{(t)}\{\pmb{x}\}, \ \mathsf{and} \ \llbracket \pmb{\delta}' \rrbracket = \left[\!\!\begin{smallmatrix} \phi \\ s \to t \end{smallmatrix} \right] \pmb{\delta}'$$

Base case: ($\delta^{(s)}$ is a constant or builtin): This case is exactly the same as in the proof of Lemma 1.

Inductive case: There are three cases to consider:

Case 1 ($\delta^{(s)}$ is a product $^{y_1} \circ \cdots \circ ^{y_K}$): This case is exactly the same as in the proof of Lemma 1.

Case 2 (top rule of $\delta^{(s)}$ is in $\mathfrak{p}^{(t)}$): This case is exactly the same as in the proof of Lemma 1.

Case 3 (top rule of $\delta^{(s)}$ is not in $\mathfrak{p}^{(t)}$): The derivation $\delta^{(s)}$ must have the following form with top rule $\mathbf{S}^{(\ell)} \in \mathcal{S}$:

$$oldsymbol{\delta}^{(s)} = egin{pmatrix} x & & & & \\ & & \mu & eta & & \\ & & & igwedge & & igwedge & & \\ & & & igwedge & & igwedge & & \\ & & & igwedge & & igwedge & & \\ & & & & & igwedge & & \\ & & & & & igwedge & & \\ & & & & & igwedge & & \\ & & & & & & igwedge & & \\ & & & & & & igwedge & & \\ & & & & & & igwedge & & \\ & & & & & & igwedge & & \\ & & & & & & igwedge & & \\ & & & & & & igwedge & & \\ & & & & & & igwedge & & \\ & & & & & & igwedge & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & \\ & & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & \\ & &$$

We now fold the derivation $\stackrel{\mu}{\bigtriangleup}$: $u^{-1}\left(\stackrel{\mu}{\bigtriangleup},\mathbf{S}^{(\ell)}\right)=\stackrel{y_j}{\bigtriangleup}$. Observe that $\stackrel{\alpha}{\smile},\stackrel{y_j}{\smile},\stackrel{\beta}{\smile} \prec$ steps $\boldsymbol{\delta}^{(s)}$ because they are subroutines.

C₂: Lemma 3 guarantees that the remapped tree $\overset{y_j}{\triangle} \in \mathfrak{p}^{(s)}\{y_j\}$. Furthermore, the folder rule r is a (distinct) rule in $\mathfrak{p}^{(t)}$ of the form $x \oplus = \alpha \circ y_j \circ \beta$. Now, since $\phi \left(\overset{\alpha}{\triangle}\right), \ \phi \left(\overset{y_j}{\triangle}\right), \ \phi \left(\overset{\beta}{\triangle}\right) \in \mathfrak{p}^{(t)}\{x\}$ by [H], it follows that $\delta^{(t)} \in \mathfrak{p}^{(t)}\{x\}$.

C₃:

$$\begin{bmatrix} \phi & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & &$$

Lemma 6. Suppose

 $\boxed{ \textbf{A_3} \quad \phi_{s \to t} : \mathfrak{p}^{(s)}\{\cdot\} \to \mathfrak{p}^{(t)}\{\cdot\} \text{ is the partial function specified in Algorithm 25.} }$

Then, $\forall \pmb{\delta}^{(t)} \in \mathfrak{p}^{(t)}\{\cdot\} : \exists ! \, \pmb{\delta}^{(s)} \in \mathfrak{p}^{(s)}\{\cdot\} : \underset{s \to t}{\phi}\left(\pmb{\delta}^{(s)}\right) = \pmb{\delta}^{(t)}$. Furthermore, the function $\underset{t \to s}{\phi}$ specified in Algorithm 32 computes $\pmb{\delta}^{(s)} = \underset{t \to s}{\phi}\left(\pmb{\delta}^{(t)}\right)$.

Proof. Our proof uses induction on \prec over $\mathfrak{p}^{(t)}\{\cdot\}$. All but one of the cases in this proof are identical to Lemma 4. Fix an arbitrary $\boldsymbol{\delta}^{(t)} \in \mathfrak{p}^{(t)}\{\cdot\}$.

Inductive hypothesis:

Base case: ($\delta^{(t)}$ is a constant or builtin): Same as the proof of Lemma 4.

Inductive case: As one might expect, the construction of $\phi_{t \to s}$ essentially inverts each step of the construction of $\phi_{s \to t}$. Thus, there are three analogous cases to consider as in Lemma 5:

Case 1 ($\delta^{(t)}$ is a product $\overset{y_1}{\triangle} \circ \cdots \circ \overset{y_K}{\triangle}$): Same as the proof of Lemma 4.

Case 2 (top rule of $\delta^{(t)} \in \mathfrak{p}^{(s)}$): Same as the proof of Lemma 4.

Case 3 (top rule of $\delta^{(t)} \notin \mathfrak{p}^{(s)}$):

We seek to identify all $\delta^{(s)}$ that satisfy the following equation:

$$oldsymbol{\delta}^{(t)} = \overbrace{ egin{array}{c} x \\ y_j \\ A \end{array} }^{oldsymbol{x}} = \phi _{s o t} \Big(oldsymbol{\delta}^{(s)} \Big)$$

It is clear from the implementation of $\phi_{s \to t}$ (see line 25.7) that $\pmb{\delta}^{(t)}$ has this form is if and only if

By IH, we have unique solutions to the following equations:

Furthermore, those solutions are given by $\phi_{t \to s}$:

$$\phi \left(\bigwedge_{t \to s}^{\alpha} \right) = \pmb{\delta}_0^{(s)} \quad ; \quad \phi \left(\bigwedge_{t \to s}^{y_j} \right) = \pmb{\delta}_1^{(s)} \quad ; \quad \phi \left(\bigwedge_{t \to s}^{\beta} \right) = \pmb{\delta}_2^{(s)}$$

Next, we solve $\delta_1^{(s)}=u^{-1}\Big(\delta_3^{(s)}\Big)$. According to Lemma 3, this equation has a unique solution that is computed by $\delta_3^{(s)}=u^{-1}\Big(\delta_1^{(s)}\Big)$ with line 25.7. Lastly, we have that the unique solution $\delta^{(s)}$ is given by

$$\boldsymbol{\delta}^{(s)} = \phi \begin{pmatrix} x \\ \alpha & y_j & \beta \\ \triangle & \triangle & \triangle \end{pmatrix} \stackrel{\text{def}}{=} \begin{pmatrix} x \\ \phi \\ \triangle & A \end{pmatrix} \quad u\begin{pmatrix} \phi \\ \psi_j \\ \psi_j \\ \psi_j \end{pmatrix} \quad \phi \begin{pmatrix} \beta \\ \triangle \end{pmatrix}$$

Theorem 3. Given $\mathfrak{p}^{(s)} \xrightarrow{\text{fold}_{S \to r}^{d,j}} \mathfrak{p}^{(t)}$ and $\mathfrak{p}^{(s)} \overset{\text{strong}}{\equiv} \mathfrak{p}^{(d)}$. Suppose the function $\phi_{s \to t} (\delta^{(s)})$, as specified in Algorithm 25, terminates for all $\delta^{(s)} \in \mathfrak{p}^{(s)}\{\cdot\}$. Then, $\phi_{s \to t}$ is a value-preserving bijection between $\mathfrak{p}^{(s)}\{\cdot\}$ and $\mathfrak{p}^{(t)}\{\cdot\}$. Furthermore, $\mathfrak{p}^{(s)} \overset{\text{strong}}{\equiv} \mathfrak{p}^{(t)}$.

Proof. Under the above conditions, it follows from Lemma 5 and Lemma 6 that ϕ is a value-preserving bijection between $\mathfrak{p}^{(s)}\{\cdot\}$ and $\mathfrak{p}^{(t)}\{\cdot\}$. Such a bijection ensures that $\mathfrak{p}^{(s)} \stackrel{\text{strong}}{\equiv} \mathfrak{p}^{(t)}$ by Definition 31.

Theorem 4. Suppose $\mathfrak{p}^{(s)} \xrightarrow{\text{unfold}} \mathfrak{p}^{(t)}$ according to Definition 33, then $\mathfrak{p}^{(s)} \stackrel{\text{strong}}{=} \mathfrak{p}^{(t)}$.

Proof. Any basic unfold transformation is an instance of a generalized unfold transformation with $\mathfrak{p}^{(d)}=\mathfrak{p}^{(s)}$. Thus, it inherits the correctness proof for the generalized unfold transformation (Theorem 2). We can strengthen the claim: all basic unfold transformations are semantics-preserving because we can drop Theorem 2's stipulation that the derivation mapping (Algorithm 24) terminates as it is always true for basic unfold transformations. To see why, we observe that the call to u always returns a (strict) subtree of the $\delta^{(s)}$, and subtrees are a well-founded order. Thus, ϕ is guaranteed to terminate for any basic unfold transformation.

Theorem 5. Suppose $\mathfrak{p}^{(s)} \xrightarrow{\mathsf{fold}} \mathfrak{p}^{(t)}$ according to Definition 34, then $\mathfrak{p}^{(s)} \stackrel{\mathsf{strong}}{=} \mathfrak{p}^{(t)}$.

Proof. Theorem 4 establishes a symmetrical result, which we can leverage here. Definition 34 requires reversibility: $\mathfrak{p}^{(t)} \xrightarrow{\mathrm{unfold}_{r \to \mathcal{S}}^{t,j}} \mathfrak{p}^{(s)}$. Since $\mathfrak{p}^{(t)} \xrightarrow{\mathrm{unfold}_{r \to \mathcal{S}}^{t,j}} \mathfrak{p}^{(s)}$ implies $\mathfrak{p}^{(t)} \stackrel{\mathrm{strong}}{\equiv} \mathfrak{p}^{(s)}$ (by Theorem 4) and semantic equivalence is symmetric, it follows that $\mathrm{fold}_{\mathcal{S} \to r}^{s,j}$ is also semantics-preserving because any transformation from $\mathfrak{p}^{(s)}$ to $\mathfrak{p}^{(t)}$ (or vice versa) would be semantics-preserving under the condition that a strongly equivalent transformation relates $\mathfrak{p}^{(t)}$ and $\mathfrak{p}^{(s)}$.

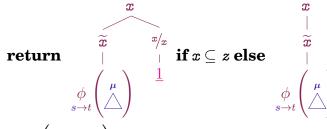
B.3 Speculation

We provide a derivation mapping as an (incomplete) justification for the correctness of the speculation transformation (Definition 41). To complete the proof of correctness, we would need to prove that this derivation mapping is a bijective value-preserving mapping (as we did in our proofs of correctness for fold and unfold in §B.2). This should be a fairly straightforward argument to make by structural induction (on the subtree relation); however, we do not provide such an argument.

Algorithm 33 Derivation mapping for $\mathfrak{p}^{(s)} \xrightarrow{\text{speculation}_{\mathbf{z},\sigma}} \mathfrak{p}^{(t)}$.

- 1. $\operatorname{def} \phi_{s \to t} \left(\delta^{(s)} \right)$:
- case $\delta^{(s)}$ is a builtin or constant: **return** $\delta^{(s)}$
- 4. **case** $\begin{pmatrix} x \\ \mu \\ \triangle \end{pmatrix}$ **and** the top rule $\notin \sigma$:

 \triangleright Bottom of the spine



- 5. **case** $\begin{pmatrix} x \\ y & \beta \\ \triangle & \triangle \end{pmatrix}$ **and** top rule $\in \sigma$:
- 6. **match** $\phi \left(\bigwedge^{y} \right)$
- 7. **case** $y : \mathbf{return}$ $\widetilde{x} : \widetilde{x}$ $\widetilde{x} : \widetilde{x} : \widetilde{$

Appendix C

Pattern Matching Algorithms

C.1 Unification

In this section, we describe algorithms for unification, which we briefly described in §3.1.1.²²³ Much has been written about unification (Herbrand, 1930; Robinson, 1965; Knight, 1989; Norvig, 1991; Martelli and Montanari, 1982) so we will not talk too much about the algorithm—we have primarily included it for completeness. Our implementation of explicit unification unify is given in Algorithm 34; it closely follows the pseudocode from the textbook of Russell and Norvig (2020). Overall, the algorithm is a relatively straightforward recursive procedure that simultaneously traverses a pair of terms from the top down, acquiring the minimal set of variable bindings needed to make the terms equal (or failing if equality is impossible). Algorithm 35 provides a "one-

 $^{^{223}}$ The implementation in this section uses explicit substitution maps, whereas in §3.1.2.1 the map is represented implicitly using a mutable data structure.

sided" unification algorithm (cover) which solves unification problems of the form $\theta(x)=y$. This algorithm can be used to determine if the set of terms denoted by x is a superset of those denoted by y: $(x\supseteq y)\iff (\mathsf{cover}(x,y)\neq \mathsf{FAIL})$. In the pseudocode, $\mathsf{deref}^*(x,\theta)$ recursively applies $\theta(x)$ until $x\notin\theta$.

Algorithm 34 unify: Unification algorithm. Solves any equation of the form $\theta(x) = \theta(y)$ for θ given possibly nonground terms x and y.

```
1. def unify(\boldsymbol{x}, \boldsymbol{y}, \boldsymbol{\theta} = \{\}):
       if \theta = \text{FAIL}: return FAIL
       else if x = y: return \theta
       else if x is a variable:
          if x \in \theta: return unify(\theta[x], y, \theta)
          else if y \in \theta: return unify(x, \theta[y], \theta)
          else if x occurs in y after substitution \theta: return FAIL
          else
              \theta[x] \leftarrow y
                                                                                                                \triangleright Mutates \theta
              return \theta
10
       else if y is a variable: return unify(y, x, \theta)
11
       else if x and y are terms with the same name and arity:
          for i \in 1 \dots |\mathbf{x}| : \boldsymbol{\theta} \leftarrow \text{unify}(\mathbf{x}_i, \mathbf{y}_i, \boldsymbol{\theta})
13.
          return \theta
14
       else
          return FAIL
16.
```

Algorithm 35 cover: One-sided unification algorithm. Solves any equation of the form $\theta(x) = y$ for θ given possibly nonground terms x and y.

```
1. \operatorname{def} \operatorname{cover}(x,y,\theta=\{\}):
2. x \leftarrow \operatorname{deref}^*(x,\theta); y \leftarrow \operatorname{deref}^*(y,\theta)
3. \operatorname{if} x = y : \operatorname{return} \theta
4. \operatorname{else} \operatorname{if} x \operatorname{is} \operatorname{a} \operatorname{variable}:
5. \operatorname{if} \theta[x] = y : \operatorname{return} \theta
7. \operatorname{else} \operatorname{if} x \operatorname{does} \operatorname{not} \operatorname{occur} \operatorname{in} \theta(y):
8. \theta[x] \leftarrow y
9. \operatorname{return} \theta
10. \operatorname{return} \operatorname{FAIL}
11. \operatorname{else} \operatorname{if} x \operatorname{and} y \operatorname{are} \operatorname{terms} \operatorname{with} \operatorname{the} \operatorname{same} \operatorname{name} \operatorname{and} \operatorname{arity}:
12. \operatorname{for} i \in 1 \dots |x| : \theta \leftarrow \operatorname{cover}(x_i, y_i, \theta)
13. \operatorname{return} \theta
14. \operatorname{return} \theta
15. \operatorname{return} \theta
16. \operatorname{return} \operatorname{FAIL} \operatorname{else} \operatorname{if} x \operatorname{and} y \operatorname{are} \operatorname{terms} \operatorname{with} \operatorname{the} \operatorname{same} \operatorname{name} \operatorname{and} \operatorname{arity}:
17. \operatorname{return} \theta \mapsto \operatorname{cover}(x_i, y_i, \theta) \mapsto \operatorname{return} \theta
```

C.2 Sub-Product Enumeration

C.2.1 Commutative Sub-Products

Algorithm 36 is a procedure to enumerate alignments between a collection of subgoals $x_1 \circ \cdots \circ x_N$ to another collection of subgoals $y_1 \circ \cdots \circ y_M$ assuming a \circ is commutative. The main usage is during folding where $x_1 \circ \cdots \circ x_N$ is the right-hand side of a folder rule and $y_1 \circ \cdots \circ y_M$ is the right-hand side of a rule that is folded.

For example,

• covers_product(f(X) * f(Y), f(A) * f(B) * f(C)) has six alignments:

$$\begin{split} &\left\langle \left\{ \mathsf{X} \mapsto \mathsf{A}, \mathsf{Y} \mapsto \mathsf{B} \right\}, \left\{ 0 \mapsto 0, 1 \mapsto 1 \right\} \right\rangle \\ &\left\langle \left\{ \mathsf{X} \mapsto \mathsf{A}, \mathsf{Y} \mapsto \mathsf{C} \right\}, \left\{ 0 \mapsto 0, 1 \mapsto 2 \right\} \right\rangle \\ &\left\langle \left\{ \mathsf{X} \mapsto \mathsf{B}, \mathsf{Y} \mapsto \mathsf{A} \right\}, \left\{ 0 \mapsto 1, 1 \mapsto 0 \right\} \right\rangle \\ &\left\langle \left\{ \mathsf{X} \mapsto \mathsf{B}, \mathsf{Y} \mapsto \mathsf{C} \right\}, \left\{ 0 \mapsto 1, 1 \mapsto 2 \right\} \right\rangle \\ &\left\langle \left\{ \mathsf{X} \mapsto \mathsf{C}, \mathsf{Y} \mapsto \mathsf{A} \right\}, \left\{ 0 \mapsto 2, 1 \mapsto 0 \right\} \right\rangle \\ &\left\langle \left\{ \mathsf{X} \mapsto \mathsf{C}, \mathsf{Y} \mapsto \mathsf{B} \right\}, \left\{ 0 \mapsto 2, 1 \mapsto 1 \right\} \right\rangle \end{split}$$

• covers_product(f(X,Y), g(Y,Z),f(I,J), g(J,I)) has one alignment: $\langle \{X \mapsto I, Y \mapsto J, Z \mapsto K\}, \{0 \mapsto 0, 1 \mapsto 1\} \rangle$

Each iterate is a pair of a substitution map θ and index-to-index alignment σ . Substitutions maps θ and alignment σ are such that

- 1. Each alignment σ is a function $[1:N] \to [1:M]$; in other words, from the indices of x to the indices of y.
- 2. $\theta(x_i) = y_{\sigma(i)}$ for $i \in [1:N]$; like in cover, θ may only bind variables in x.

Algorithm 36 covers_product: Lazily enumerate commutative matches from a collection of subgoals $x_1 \circ \cdots \circ x_N$ to another collection of subgoals $y_1 \circ \cdots \circ y_M$. Each iterate is a pair of a substitution mapping and an index-to-index alignment.

```
1. \mathbf{def} \  \, \mathrm{covers\_product}(x_1 \circ \cdots \circ x_N, y_1 \circ \cdots \circ y_M, \theta_0 = \{\}) \colon \\ \mathbf{def} \  \, f(i, \beta, \theta) \colon \\ \mathbf{if} \  \, \theta = \mathrm{FAIL} \colon \mathbf{return} \qquad \qquad \triangleright \, \mathit{Failure} \\ \mathbf{else} \  \, \mathbf{if} \  \, |x_{[i:]}| > |y_{\beta}| \colon \mathbf{return} \qquad \qquad \triangleright \, \mathit{Failure} \colon \mathit{too} \, \mathit{many} \, \mathit{subgoals} \\ \mathbf{5.} \qquad \mathbf{else} \  \, \mathbf{if} \  \, |x_{[i:]}| = 0 \colon \mathbf{yield} \  \, \langle \theta, \{\} \rangle \qquad \qquad \triangleright \, \mathit{Success:} \, \mathit{reached} \, \mathit{the} \, \mathit{end} \, \mathit{of} \, x'\mathit{s} \\ \mathbf{6.} \qquad \mathbf{else} \\ \mathbf{7.} \qquad \mathbf{for} \  \, j \in \beta \colon \\ \mathbf{8.} \qquad \mathbf{for} \  \, \langle \theta', \sigma' \rangle \in f(i+1, \beta \smallsetminus \{j\} \,, \mathsf{cover}(x_i, y_j, \mathsf{copy}(\theta))) \colon \\ \mathbf{9.} \qquad \mathbf{yield} \  \, \langle \theta', \sigma' \cup (i \to j) \rangle \\ \mathbf{10.} \qquad \mathbf{yield} \  \, \mathbf{from} \, f(1, \{1, ..., M\} \,, \theta_0) \\ \end{cases}
```

C.2.2 Noncommutative Sub-Products

Algorithm 37 is a procedure to enumerate subproducts; it will find all the subproducts $x_1 \circ \cdots \circ x_N$ contained in $y_1 \circ \cdots \circ y_M$. Unlike the commutative case, no reordering or gaps in the match are allowed.

Algorithm 37 covers_product: Lazily enumerate non-commutative matches from a collection of subgoals $x_1 \circ \cdots \circ x_N$ to another collection of subgoals $y_1 \circ \cdots \circ y_M$. Each iterate is a pair of a substitution mapping and the start position of the match.

```
1. \mathbf{def} covers_product(x_1 \circ \cdots \circ x_N, y_1 \circ \cdots \circ y_M, \boldsymbol{\theta} = \{\}):
2. \mathbf{for} \ j \in [1:M-N]:
3. \boldsymbol{\theta}' \leftarrow \operatorname{cover}(x, y_{[j+i:j+N]}, \operatorname{copy}(\boldsymbol{\theta}))
4. \mathbf{if} \ \boldsymbol{\theta}' \neq \operatorname{FAIL}: \mathbf{yield} \ \langle \boldsymbol{\theta}', j \rangle
```

Appendix D

Dyna 2

We give a brief overview of the Dyna 2 language and thoughts about how to extend this dissertation's results to it.

Dyna 2 (Eisner and Filardo, 2011) extends Dyna 1 (cf. § 2) by allowing values to be combined in more than the two ways that a single semiring allows. Thus, Dyna 2 program can be used as a concise notation for specifying general computation graphs²²⁴ via rules that are still written in a Prolog-like notation (as with Dyna 1). Each rule continues to be a template whose *variables* may be instantiated with arbitrary terms. The rule instantiations collectively define the *contributions* to the *values* of many *items*. The contributions are *aggregated* as in Dyna 1 with an *aggregation operator* (*aggregator*, for short). Dyna programs support general mathematical expressions (not just semiring operations). A feed-forward neural network and its training objective can be defined via

²²⁴Other notations and libraries for defining computation graphs include Theano (Theano Development Team, 2016), PyTorch (Paszke et al., 2017), TensorFlow (Abadi et al., 2016), JAX (Bradbury et al., 2018), and Named Tensor Notation (Chiang, Rush, et al., 2021).

```
\begin{array}{lll} & \sigma(X) = 1/(1 + \exp(-X)) \,. & \text{$\%$ define sigmoid function} \\ & \text{$685$ out(J) = $\sigma(in(J)) \,.} & \text{$\%$ apply sigmoid function} \\ & \text{$686$ in(J) += out(I) * edge(I,J) .} & \text{$\%$ vector-matrix product} \\ & \text{$1085$ += $(out(J) - target(J))^2 .} & \text{$\%$ L2 loss} \\ \end{array}
```

The structure and weights of the network are specified by defining the values of items of the form edge(I,J). This could be done by listing one rule per edge, but it can also be done systematically by writing edge-defining rules in terms of structured node names, where these names will instantiate I,J above. For example,

```
edge(input(X,Y),hidden(X+DX,Y+DY)) = weight(DX,DY).

weight(DX,DY) := random(\star,-1,1) for DX:-4..4, DY:-4..4.
```

defines a convolutional layer with initially random weights²²⁵ and a 9×9 convolution filter. For example, the node named hidden(10,10) has a connection from the node named input(8,11) with weight weight(-2,1). The input to the network is specified at runtime by updating the values of out(input(...)); these items otherwise have no value as the input nodes have no incoming edges.

The Dyna solver seeks a *fixed point* in which all items are consistent, i.e., each item's value matches its definition from other values. A cyclic program may not have a unique fixed point. The solver is permitted to choose arbitrarily if there is more than one. In general, the solver may fail to terminate because no fixed point exists or because it cannot discover one. Sometimes the solver terminates only at numerical convergence, as for the geometric series sum

```
690 a += 1.
691 a += a/2.
```

Or this program that implements the Babylonian method for numerically

²²⁵Each expression random(*,-1,1) names a distinct random variate, since the special argument * generates a symbol ("gensym") that is different in each instantiation of the rule. Thus, the last line defines 81 distinct weights.

approximating square root: $y = \sqrt{x}$:

```
692 y := 1. % arbitrary initial value \neq 0
693 y := 1/2 * (x/y + y).
```

Here the := aggregator is the "last rule wins" aggregator. The top rule specifies the initialization for fixpoint iteration.

Eisner and Filardo (2011) give a more complete introduction to Dyna 2 with more complex examples, including message passing in variational inference and graphical models, constraint propagation in backtracking search, backpropagation, and Markov chain Monte Carlo.

Extending this Dissertation to Dyna 2

Execution. The naive, semi-naive, and prioritized fixpoint iterations algorithms in §5 have analogs in the Dyna 2 setting. They have been studied in Filardo (2017) and to some extent in Eisner and Blatz (2007).²²⁶ In settings where there is a semiring subprogram, any of our execution algorithms could be used to solve it. That said, the Dyna 2 framework appears to be moving in the direction of term-rewriting systems as the basis of its execution system (Francis-Landau et al., 2020).

Type Analysis. Our type analyzer §6.1 can be extended to Dyna 2. The main difference is that we would need to prove a type for the *value* that an item may have since it is no longer homogeneous. In other words, we would

²²⁶Eisner and Blatz (2007) described forward chaining for mixed semirings, which covers most but not all Dyna 2 programs.

reason about possible value—item *pairs*. In principle, this should not be a difficult extension. Our approach to space and time complexity analysis method (§6.3) should also work in the case of Dyna 2, as the complexity of executing a Dyna 2 program is also proportional to the cost of instantiating its rules. Our program specialization methods (§6.2) should also work with little modification.

Program Transformations. Although Eisner and Blatz (2007) predate Dyna 2, they had some foresight in their approach. Essentially, they devised semiring transformations that could be applied in mixed-semiring programs. Their trick was simply to be careful to check that there is a *locally* distributive property in a program fragment. Essentially our transformations should still work, but we just need to check for local distributivity. We also note that Dyna 2 programs will require other kinds of transformations that pull from a wider range of mathematical rewrites. Such transformations are used in other computation graph libraries.

Program Optimization by Transform Search. Our transformation search approach is directly applicable to a transformation search space tailored to Dyna 2 programs.

Appendix E

Programs Used for Experiments

E.1	Chain	s (chain-05, chain-10, bad-chain-05, bad-chain-10)	275
E.2	Projec	tive Dependency Parsing	276
	E.2.1	Bilexical Unlabeled (bilexical-unlabeled)	276
	E.2.2	Bilexical Labeled (bilexical-labeled)	277
	E.2.3	Arc Eager (arc-eager)	277
	E.2.4	Split Bilexical Dependency Parsing	278
		E.2.4.1 Johnson (2007)'s Version (split-head-J)	278
		E.2.4.2 Eisner and Blatz (2007)'s Version (split-head-EB)	278
E.3	Const	ituency Parsing	279
	E.3.1	CKY (cky3)	279
	E.3.2	CKY with 4-ary Productions (cky4)	279
	E.3.3	CKY with a Fixed Grammar (cky+grammar)	280
	E.3.4	Tensor Decomposition Parser (tensor-decomp-parser)	281
E.4	Invers	sion Transduction Grammars (itg)	281

E.5	Edit Distance (edit)	282
E.6	Bar-Hillel Construction (bar-hillel)	282
E.7	Expectations under a Linear-Chain Conditional Random Field	
	(chain-expect)	283
E.8	Intersection of Even and Odd Peano Numbers (even-odd-peano) .	284
E.9	Hidden Markov Models (hmm)	285
E.10	Semi-Markov Model (semi-markov)	285
E.11	Shortest Path in a Graph	286
	E.11.1 Explicit Path (path-list)	286
	E.11.2 Start Path (path-start)	286
E.12	Pushdown Automaton with an Explicit Stack (explicit-pda)	287

E.1 Chains (chain-05, chain-10, bad-chain-05, bad-chain-10)

Chain structures are common in NLP as they model the interaction between adjacent words or labels in a sequence (e.g., Lafferty et al. (2001)).

The chain-N programs are adaptations of Example 3,

```
goal += w(X_1, X_2) * w(X_2, X_3) * w(X_3, X_4) * w(X_4, X_5).
```

which is for the specific case of length 4, to have length N. The original program's degree is N. The optimal degree for chains is 2 regardless of N.

We also have the following examples of bad variable-elimination orderings

§7.3. Both of the examples below can be sped up to degree 2.

bad-chain-05:

```
1 goal += gen6(X5) * tmp(X1,X4,X5).

2 tmp(X1,X4,X5) += gen2(X1,X4) * w(X4, X5).

3 gen2(X1,X4) += gen5(X1,X2) * gen4(X2, X4).

4 gen4(X2,X4) += w(X2, X3) * w(X3, X4).

5 gen5(X1,X2) += w(X1, X2).

6 gen6(X5) += w(X5, X6).
```

bad-chain-10:

```
1 gen56(X_4, X_6) += w(X_4, X_5) * w(X_5, X_6).

2 gen57(X_2) += w(X_1, X_2).

3 gen59(X_4, X_7) += gen56(X_4, X_6) * w(X_6, X_7).

4 goal += gen124(X_3) * gen122(X_3).

5 gen124(X_3) += gen123(X_4) * w(X_3, X_4).

6 gen123(X_4) += gen125(X_{10}) * gen121(X_{10}, X_4).

7 gen125(X_{10}) += w(X_{10}, X_{11}).

8 gen122(X_3) += gen57(X_2) * w(X_2, X_3).

9 gen121(X_{10}, X_4) += gen119(X_4, X_8) * gen120(X_{10}, X_8).

10 gen120(X_{10}, X_8) += w(X_8, X_9) * w(X_9, X_{10}).

11 gen119(X_4, X_8) += gen59(X_4, X_7) * w(X_7, X_8).
```

E.2 Projective Dependency Parsing

This program performs a version of CKY (§E.3) where the non-terminals have been lexically annotated with their head word. In contrast to CKY's $\mathcal{O}(n^3)$ runtime, the original presentation of this algorithm ran in $\mathcal{O}(n^5)$ (Collins, 1996). In subsequent work, Eisner and Satta (1999a) gave a faster version of the algorithm that runs in $\mathcal{O}(n^4)$.

E.2.1 Bilexical Unlabeled (bilexical-unlabeled)

```
phrase(I,H,K) += phrase(I,H,J) * phrase(J,H',K) * score(H \leftarrow H').
phrase(I,H',K) += phrase(I,H,J) * phrase(J,H',K) * score(H \rightarrow H').
phrase(I,I,K) += word(I,_K).
```

```
4 goal(N) += phrase(0,0,N).

5 inputs: word(_,_,_); score(_ \rightarrow _); score(_ \leftarrow _).

7 outputs: goal(_).
```

Degree: 5. Optimal: 4.

E.2.2 Bilexical Labeled (bilexical-labeled)

Extends bilexical-unlabeled with labels (i.e., grammar relations).

```
phrase(I,H,X,K) += phrase(I,H,Y,J) * phrase(J,H',Z,K) * rewrite(H,Y, H',Z, H,X).
phrase(I,H',X,K) += phrase(I,H,Y,J) * phrase(J,H',Z,K) * rewrite(H,Y, H',Z, H',X).
phrase(I,I,X,K) += rewrite(X,W) * word(I,W,K).
goal(N) += phrase(0,0,s,N).

inputs: word(_,_,_); rewrite(_,_,_,_,_); rewrite(_,_).
outputs: goal(_).
```

Degree: 8. Optimal: 7.

E.2.3 Arc Eager (arc-eager)

This benchmark is based on Example 21; please refer to that discussion.

```
1 \beta(I,I^0,K) += word(I,K).

2 \beta(I,I^1,K) += word(I,K).

3 \beta(I,H^0,K) += \beta(I,H^0,J) * score(H \leftarrow H') * \beta(J,H'^0,K).

4 \beta(I,H^0,K) += \beta(I,H^0,J) * score(H \rightarrow H') * \beta(J,H'^1,K).

5 goal(N) += \beta(0,0^0,N).

6 r inputs: word(_,_); score(_ \rightarrow _); score(_ \leftarrow _).

8 outputs: goal(_).
```

Degree: 6, Optimal: 3.

E.2.4 Split Bilexical Dependency Parsing

A split bilexical grammar is one where the headword must combine with all its right children before any of its left children. Johnson (2007) and Eisner and Satta (1999a) each gave initial, inefficient algorithms that they show how to speed up to $\mathcal{O}(n^3)$. Their initial algorithms are different, so, we have provided them as separate benchmarks.

E.2.4.1 Johnson (2007)'s Version (split-head-J)

The following program is the initial, inefficient algorithm of Johnson (2007).

```
1 goal(N) += x(0,V,N).
2
3 % words are duplicated, as in Johnson (2007).
4 l(I,K) += word(left,I,K).
5 r(I,K) += word(right,I,K).
6
7 % l(I, K) is a span with K as the head
8 l(I,K) += x(I,V,J) * score(V ← K) * l(J,K).
9
10 % r(I, K) is a span with I as the head
11 r(I,K) += r(I,J) * score(I → V) * x(J,V,K).
12
13 % J is the head of l(I, J) and J is the head of r(J, K)
14 x(I,J,K) += l(I,J) * r(J,K).
15
16 input: word(_,_,_); score(_ → _); score(_ ← _).
17 output: goal(_).
```

Degree: 4, Optimal: 3.

E.2.4.2 Eisner and Blatz (2007)'s Version (split-head-EB)

The following program is the initial, inefficient algorithm of Eisner and Blatz (2007), which we discussed in Example 44.

```
1 r(I,I,K) += word(I,K). % no right children yet
```

Degree: 5, Optimal: 3.

E.3 Constituency Parsing

The following programs are variants of CKY (Cocke and Schwartz, 1970; Younger, 1967; Kasami, 1965; Lange and Leiß, 2009; Eisner and Blatz, 2007; Tomita, 1985; Tomita, 1991; Johnson, 1989; Baker, 1979; Jelinek, 1985). We briefly discussed CKY in Example 6 of the main text.

E.3.1 CKY (cky3)

```
1 \beta(X,I,K) += \gamma(X,Y,Z) * \beta(Y,I,J) * phrase(Z,J,K).

2 \beta(X,I,K) += \gamma(X,Y) * \beta(Y,I,K).

3 \beta(X,I,K) += \gamma(X,Y) * word(Y,I,K).

4 z += \beta(root, 0, N) * len(N).

5 input: word(W,I,K); len(N); \gamma(X,Y,Z); \gamma(X,Y).

6 output: z.
```

Degree: 6, Optimal: 5.

E.3.2 CKY with 4-ary Productions (cky4)

The following program implements Tomita (1985)'s algorithm where the grammar can have a production rule of length up to 4. It is a simple modification

to CKY4. We just add the following rule and input declaration.

```
7 phrase(X,I<sub>1</sub>,I<sub>4</sub>) += \gamma(X,Y<sub>1</sub>,Y<sub>2</sub>,Y<sub>3</sub>) * phrase(Y<sub>1</sub>,I<sub>1</sub>,I<sub>2</sub>) * phrase(Y<sub>2</sub>,I<sub>2</sub>,I<sub>3</sub>)

8 * phrase(Y<sub>3</sub>,I<sub>3</sub>,I<sub>4</sub>).

9 input: \gamma(X,Y<sub>1</sub>,Y<sub>2</sub>,Y<sub>3</sub>)
```

Degree: 7, Optimal: 6.

E.3.3 CKY with a Fixed Grammar (cky+grammar)

For CKY+grammar, we use CKY3, but we remove the input declaration for γ and declare the following grammar for the program to specialize to.

```
10 \gamma("S", "NP", "VP") += 1.0.
   \gamma("NP", "Det", "N") += 1.0.
\gamma("NP", "NP", "PP") += 1.0.
13 \gamma("VP", "V", "NP") += 1.0.
   \gamma("VP", "V") += 1.0.
   \gamma("VP", "VP", "PP") += 1.0.
\gamma("PP", "P", "NP") += 1.0.
                   "NP") += 1.0.
\gamma("NP", "Papa") += 1.0.
            "caviar") += 1.0.
   γ("N",
20 \gamma("N", "spoon") += 1.0.
\gamma("N", "fork") += 1.0.
   \gamma("N", "telescope") += 1.0.
   \gamma("N", "boy") += 1.0.
\gamma("N", "girl") += 1.0.
\gamma("N", "baby") += 1.0.
   \gamma("N", "man") += 1.0.
\gamma("N", "woman") += 1.0.
28 \gamma("N", "moon") += 1.0.
   \gamma("N", "cat") += 1.0.
30 \gamma("V", "ate") += 1.0.
31 \gamma("V", "saw") += 1.0.
\gamma("V", "fed") += 1.0.
   \gamma("V", "said") += 1.0.
   \gamma("V", "jumped") += 1.0.
   \gamma("P", "with") += 1.0.
   \gamma("P", "over") += 1.0.
   \gamma("P", "under") += 1.0.
38 \gamma("P", "above") += 1.0.
39 \gamma("P", "below") += 1.0.
40 \gamma("P", "on") += 1.0.
   \gamma("P", "in") += 1.0.
```

Degree: 6, Optimal: 3.

E.3.4 Tensor Decomposition Parser (tensor-decompparser)

The following program is a variation on CKY where the grammar (i.e., the rewrite relation) is factored: $rewrite(X,Y,Z) += \lambda(R) * u(R, X) * v(R, Y) * w(R, Z)$. This approach was proposed by Cohen, Satta, et al. (2013) whereby the grammar was approximated using low-rank tensor decomposition and subsequently parsed with the algorithm below. However, the algorithm below can be sped up by exploiting the structure of the rewrite relation.

Degree: 6, Optimal: 4.

E.4 Inversion Transduction Grammars (itg)

Inversion transduction grammars were introduced by Wu (1996), who gave an $\mathcal{O}(n^7)$ algorithm, which was later sped up to $\mathcal{O}(n^6)$ by Huang, Zhang, et al. (2005). The model simultaneously parses a pair of related sentences—typically a target language sentence and a source language sentence, as in machine translation. The model allows for a restricted syntactic reordering of phrases

called inversion.

```
\begin{array}{lll} & \beta(A,I_1,K_1,I_2,K_2) \ += \ word_1(X,I_1,K_1) \ * \ word_2(X',I_2,K_2) \ * \ transduce(A,X,X'). \\ & \beta(A,I_1,K_1,I_2,K_2) \ += \ \beta(B,I_1,J_1,I_2,J_2) \ * \ \beta(C,J_1,K_1,J_2,K_2) \ * \ rewrites(A,B,C). \\ & \beta(A,I_1,K_1,I_2,K_2) \ += \ \beta(B,J_1,K_1,I_2,J_2) \ * \ \beta(C,I_1,J_1,J_2,K_2) \ * \ rewritesinv(A,B,C). \\ & goal \ += \ \beta(s,0,N_1,0,N_2) \ * \ length(N_1,N_2). \\ & inputs: \ word_1(\_,\_,\_); \ word_2(\_,\_,\_); \ transduce(\_,\_,\_); \ rewritesinv(\_,\_,\_); \\ & rewritesinv(\_,\_,\_); \ length(\_,\_). \\ & output: \ goal. \\ \end{array}
```

Degree: 9, Optimal: 8.

E.5 Edit Distance (edit)

The following program implements a weighted generalized monotonic alignment between two sequences word and word. It is essentially the well-known Levenstein distance (Levenshtein, 1966).

```
align(0,0) min= 1.
align(J,J') min= align(I,I') * word(W,I,J) * word'(W',I',J') * score(W,W').
align(I,J') min= align(I,I') * word'(W',I',J') * score(ε,W').
align(J,I') min= align(I,I') * word(W,I,J) * score(W,ε).
goal min= align(N,N') * len(N) * len'(N').
input: word(_,_,_); word'(_,_,_); score(_,_); len(_); len'(_).
output: goal.
```

Degree 6, Optimal: 4.

E.6 Bar-Hillel Construction (bar-hillel)

The following program implements a parser for the (weighted) intersection of a context-free parser and a bigram model on the part-of-speed sequences. It is essentially Bar-Hillel et al. (1961)'s construction of a context-free language that accepts the intersection of a regular language and a context-free language.

Degree 10, Optimal: 8.

E.7 Expectations under a Linear-Chain Conditional Random Field (chain-expect)

This example implements the inside-outside speedup (Li and Eisner, 2009) for computing the expectation of an additively decomposable function $f: S \times S' \to \mathbb{R}^d$ over randomly drawn sequences from a weighted graph (e.g., a conditional random field (Lafferty et al., 2001)). The graph is specified as a collection of weights w, as well as start and end nodes. The relations α and β implement the forward-backward algorithm (discussed in Rabiner (1989)), and z is the normalization constant of the distribution. The expectation of the i^{th} dimension of f is fbar(I)/z.

```
1 % forward algorithm
2 \alpha(S) += start(S).
3 \alpha(S') += \alpha(S) * w(S,S').
4 % backward algorithm
5 \beta(S) += end(S).
6 \beta(S) += w(S,S') * \beta(S').
7 % normalization constant
8 z += \alpha(S) * end(S).
9 % unnormalized expected value via inside-outside speedup
10 fbar(R) += \alpha(S) * w(S,S') * \beta(S') * r(S,S',R).
11 input: w(-,-); r(-,-,-); start(_); end(_).
```

```
output: fbar(_). z.
```

Degree 3, Optimal: 3.

E.8 Intersection of Even and Odd Peano Numbers (even-odd-peano)

The following program counts the number of Peano numbers that are both even and odd. No such numbers exist, but as the program is written, fixpoint iteration will enumerate the even and odd numbers forever. Thus, its running time is infinite. However, it is possible to transform this program into a version that runs in constant time using fold/unfold and define; the optimal program is the empty program.

```
1  % even peano numbers
2  even(0).
3  even(s(s(X))) += even(X).
4
5  % odd peano numbers
6  odd(X) += even(s(X)).
7
8  % intersection
9  evenodd(X) += even(X) * odd(X).
10
11  % size of the intersection
12  goal += evenodd(X).
13
14  output: goal.
```

Degree: 1, Optimal: 0.

E.9 Hidden Markov Models (hmm)

Hidden Markov models (HMMs) are the generative and locally normalized analog of CRFs, which are discussed in §E.7. Rabiner (1989) provides a classic tutorial.

```
v(0,start) += 1.
v(T',Y') += v(T,Y) * emission(Y,X) * transition(Y,Y') * obs(T,X,T').
goal += v(N,stop) * len(N).
input: obs(_,_,_); len(_); emission(_,_); transition(_,_).
output: goal.
```

Degree 5, Optimal: 4.

E.10 Semi-Markov Model (semi-markov)

The semi-Markov model (Sarawagi and Cohen, 2004) generalizes a Markov model to score spans rather than individual words. In terms of runtime, one can perform inference in a Markov model in $\mathcal{O}(n)$ time (omitting dependence on the number of tags). In contrast, inference in a semi-Markov model takes $\mathcal{O}(n^2)$.

```
1 \beta(\text{start}, 0) += 1.

2 \beta(Y, J) += \beta(X, I) * \text{transition}(X, Y) * \text{chunk}(Y, I, J).

3 \text{goal} += \beta(\_, N) * \text{len}(N).

4 \text{input: transition}(\_,\_); \text{chunk}(\_,\_,\_); \text{len}(\_).

5 \text{output: goal}.
```

Degree 4, Optimal: 3.

E.11 Shortest Path in a Graph

E.11.1 Explicit Path (path-list)

The code below was previously discussed in Example 28. This example is interesting because its initial running time is infinite, but it can be sped up to $\mathcal{O}(n^2)$ where n is the number of nodes in the graph. This example requires either generalized fold/unfold or the left-corner transformation (Example 46)

```
1 β([X,Y|Xs]) += edge(X,Y) * β([Y|Xs]).
2 β([X]) += stop(X).
3 goal += start(X) * β([X|Xs]).
4
5 outputs: goal.
6 inputs: start(_); edge(_,_); stop(_).
```

E.11.2 Start Path (path-start)

The code below was previously discussed in Example 27. Here we speed up an initial $\mathcal{O}(n^3)$ program to $\mathcal{O}(n^2)$ where n is the number of nodes in the graph. This example requires generalized fold/unfold.

```
path(I,I).
path(I,K) += path(I,J) * edge(J,K).

goal += start(I) * path(I,K) * stop(K).

outputs: goal.
inputs: start(_); edge(_,_); stop(_).
```

E.12 Pushdown Automaton with an Explicit Stack (explicit-pda)

The program below can, in principle, be converted into CKY. However, doing so requires program transformations that appear to go beyond those in our arsenal.

The following program was manually specified using an extension of speculation, which permits a CCG-like composition rule.²²⁷ Unfortunately, we have not been able to fully formalize the correctness conditions for this extended transformation as of the time of writing. We can, however, use our fold/unfold transformations to prove that the program below is equivalent to the original.

Call the original program $\mathfrak{p}^{(0)}$ and the program below $\mathfrak{p}^{(1)}$.

Recovering CKY. First, let us see how the program converts into CKY. If we unfold the goal rule (twice), we arrive at a program that no longer depends

²²⁷It appears as if our hoisting transformations are incapable of this transformation because they only have a CCG-like *application* rule.

on the pda relation:

```
pda(Ys,I) / pda(Ys,I).
pda([X|Xs],K)/pda(Xs,I) += pda(Xs,J)/pda(Xs,I) * rewrite(X,Y) * word(Y,J,K).
pda([X|Xs],K)/pda(Xs,I) +=
pda([Z,Y|Xs],K)/pda([Y|Xs],J) * pda([Y|Xs],J)/pda(Xs,I) * rewrite(X,Y,Z).
goal += len(N) * pda([s],N)/pda([],0).
```

This program does not instantiate the xs. We could run this program as it is, and it will have the $\mathcal{O}(k^3 \, n^3)$ running time of CKY. We can also use our program specialization §6.2 transformation to arrive at the following program, which is visually more similar to CKY. Notice that the argument ordering to the new relation p is chosen arbitrarily.

```
 \begin{array}{lll} & p(X,K,I) \; += \; rewrite(X,Y) \; * \; word(Y,I,K) \, . \\ & p(X,K,I) \; += \; p(Z,K,J) \; * \; p(Y,J,I) \; * \; rewrite(X,Y,Z) \, . \\ & goal \; += \; len(N) \; * \; p(s,N,\emptyset) \, . \end{array}
```

It is easy to see that this program is equivalent to CKY.

Correctness by fold/unfold. Now, we investigate the correctness of our manually specified program. The key is to show that the recovery rules can be transformed into the original definition $\mathfrak{p}^{(0)}$. If we can do that using semantics-preserving transformations, we would have a proof that $\mathfrak{p}^{(0)}$ and $\mathfrak{p}^{(1)}$ are equivalent.

```
pda([X|Ys],K) += (pda([X|Ys],K)/pda(Ys,I)) * pda(Ys,I).
pda([],0).
```

Can be transformed into the original rules:

```
pda([],0).
pda([X|Xs],K) += pda(Xs,J) * rewrite(X,Y) * word(Y,J,K).
pda([X|Xs],K) += pda([Y2,Y1|Xs],K) * rewrite(X,Y1,Y2).
```

Unfolding the recursive rule twice (using $\mathfrak{p}^{(1)}$ as auxiliary definitions):

```
pda([],0).
pda([X|Ys],K) += pda([Z,Y|Ys],K) / pda([Y|Ys],J) * pda([Y|Ys],J) / pda(Ys,I) *
```

```
rewrite(X,Y,Z) * pda(Ys,I).
pda([X|Ys],K) += rewrite(X,Y) * word(Y,I,K) * pda(Ys,I).
```

Folding this program twice (using $\mathfrak{p}^{(1)}$ as auxiliary definitions) us what we were after.

```
pda([],0).
pda([X|Ys],K) += pda([Z,Y|Ys],K) * rewrite(X,Y,Z).
pda([X|Ys],K) += rewrite(X,Y) * word(Y,I,K) * pda(Ys,I).
```

Since we used semantics-preserving rules to arrive at this definition, it follows that $\mathfrak{p}^{(0)} \equiv \mathfrak{p}^{(1)}$. We used our Dyna library to work through this example, meaning we did not have to manually perform the transformations or check their safety conditions; all that was handled seamlessly by the library. We used our search algorithm to enumerate the valid folds but chose among them manually.

References

- Abadi, Martin, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, Manjunath Kudlur, Josh Levenberg, Rajat Monga, Sherry Moore, Derek G. Murray, Benoit Steiner, Paul Tucker, Vijay Vasudevan, Pete Warden, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng (2016). "TensorFlow: A System for Large-scale Machine Learning". In: USENIX Symposium on Operating Systems Design and Implementation (OSDI 16). URL: https://www.usenix.org/system/files/conference/osdi16/osdi16-abadi.pdf.
- Aït-Kaci, Hassan (1991). Warren's Abstract Machine: A Tutorial Reconstruction.

 MIT Press. URL: https://web.archive.org/web/20031004100957/http://www.vanx.org/archive/wam/wambook.pdf.
- Aji, S. M. and R. J. McEliece (2006). "The Generalized Distributive Law". In: IEEE Transactions on Information Theory 46.2. URL: https://doi.org/10.1109/ 18.825794.
- Alonso, Miguel A., David Cabrero, Eric de la Clergerie, and Manuel Vilares (1999). "Tabular Algorithms for TAG Parsing". In: *The Conference of the*

- European Chapter of the Association for Computational Linguistics. URL: https://aclanthology.org/E99-1020.
- Althoen, Steven C. and Renate McLaughlin (1987). "Gauss-Jordan Reduction:

 A Brief History". In: *The American Mathematical Monthly* 94.2. URL: http://www.jstor.org/stable/2322413.
- Apt, Krzysztof R. and Mark Wallace (2007). Constraint Logic Programming using ECLiPSe. Cambridge University Press. URL: https://doi.org/10.1017/CB09780511607400.
- Bagnara, Roberto, Fred Mesnard, Andrea Pescetti, and Enea Zaffanella (2010).

 "The Automatic Synthesis of Linear Ranking Functions: The Complete Unabridged

 Version". In: Computing Research Repository abs/1004.0944. arXiv: 1004.0944.

 URL: http://arxiv.org/abs/1004.0944.
- Baharev, Ali, Hermann Schichl, and Arnold Neumaier (2016). "Decomposition Methods for Solving Sparse Nonlinear Systems of Equations". In: *In submission*. URL: https://baharev.info/publications/baharev_tearing_survey.pdf.
- Baker, J. K. (1979). "Trainable Grammars for Speech Recognition". In: Speech Communication Papers Presented at the Meeting of the Acoustical Society of America. URL: https://compcalc.github.io/public/baker/baker1979trainable.pdf.
- Baldoni, Roberto, Emilio Coppa, Daniele Cono D'elia, Camil Demetrescu, and Irene Finocchi (2018). "A Survey of Symbolic Execution Techniques". In: *ACM Computing Surveys* 51.3. URL: https://doi.org/10.1145/3182657.

- Bar-Hillel, Yehoshua, Micha Asher Perles, and Eli Shamir (1961). "On Formal Properties of Simple Phrase-Structure Grammars". In: Zeitschrift für Phonetik, Sprachwissenschaft und Kommunikationsforschung 14 (2). URL: https://www.proquest.com/openview/fb41296047fb7453dcb1de182b4aa0b6/1.
- Baras, John S. and George Theodorakopoulos (2010). *Path Problems in Networks*. Synthesis Lectures on Communication Networks. Morgan & Claypool Publishers. URL: https://doi.org/10.2200/S00245ED1V01Y201001CNT003.
- Bellman, Richard (1958). "On a Routing Problem". In: Quarterly of Applied

 Mathematics 16.1. URL: http://www.jstor.org/stable/43634538.
- Bilmes, Jeff, Kste Asanović, Jim Demmel, Dominic Lam, and Chee-Whye Chin (1996). *PHiPAC: A Portable, High-Performance, ANSI C Coding Methodology* and its Application to Matrix Multiply. Tech. rep. University of Tennessee.

 URL: https://people.ece.uw.edu/bilmes/p/pgs/b2hd-bilmes1996-lapack.html.
- Bilmes, Jeff and Geoffrey Zweig (2002). "The Graphical Models Toolkit: An Open Source Software System for Speech and Time-Series Processing". In: 2002 IEEE International Conference on Acoustics, Speech, and Signal Processing.

 Vol. 4. URL: https://ieeexplore.ieee.org/document/5745513.
- Bossi, Annalisa, Nicoletta Cocco, and Sandro Etalle (1992). "On Safe Folding".

 In: International Symposium on Programming Language Implementation

 and Logic Programming. URL: https://link.springer.com/chapter/10.1007/3-540-55844-6_135.
- Bradbury, James, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas,

- Skye Wanderman-Milne, and Qiao Zhang (2018). *JAX: Composable Transformations of Python+NumPy Programs*. Version 0.3.13. URL: http://github.com/google/jax.
- Büchse, Matthias, Mark-Jan Nederhof, and Heiko Vogler (2011). "Tree Parsing with Synchronous Tree-Adjoining Grammars". In: *International Workshop on Parsing Technologies*. URL: https://www.aclweb.org/anthology/W11-2903.
- Burden, Håkan and Peter Ljunglöf (2005). "Parsing Linear Context-Free Rewriting Systems". In: *Proceedings of the International Workshop on Parsing Technology*. URL: https://aclanthology.org/W05-1502.
- Burstall, Rod M. and John Darlington (1977). "A Transformation System for Developing Recursive Programs". In: *Jorunal of the ACM* 24.1. URL: https://doi.org/10.1145/321992.321996.
- Chekuri, Chandra and Anand Rajaraman (1997). "Conjunctive Query Containment Revisited". In: *Proceedings of the International Conference on Database Theory*. URL: https://doi.org/10.1007/3-540-62222-5_36.
- Cheney, James, Laura Chiticariu, and Wang-Chiew Tan (2009). "Provenance in Databases: Why, How, and Where". In: Foundations and Trends® in Databases 1.4. URL: http://dx.doi.org/10.1561/1900000006.
- Chiang, David, Colin McDonald, and Chung-chieh Shan (2023). "Exact Recursive Probabilistic Programming". In: URL: https://arxiv.org/abs/2210.01206.
- Chiang, David, Alexander M. Rush, and Boaz Barak (2021). "Named Tensor Notation". In: Computing Research Repository abs/2102.13196. arXiv: 2102. 13196. URL: https://arxiv.org/abs/2102.13196.

- Cocke, John and Jacob T. Schwartz (1970). *Programming Languages and their Compilers: Preliminary Notes*. Tech. rep. Courant Institute of Mathematical Sciences, New York University. URL: http://www.softwarepreservation.org/projects/FORTRAN/CockeSchwartz_ProgLangCompilers.pdf.
- Cohen, Shay B., Carlos Gómez-Rodríguez, and Giorgio Satta (2011). "Exact Inference for Generative Probabilistic Non-Projective Dependency Parsing". In: Proceedings of the Conference on Empirical Methods in Natural Language Processing. URL: https://aclanthology.org/D11-1114.
- Cohen, Shay B., Giorgio Satta, and Michael Collins (2013). "Approximate PCFG Parsing Using Tensor Decomposition". In: *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics:*Human Language Technologies. URL: https://aclanthology.org/N13-1052.
- Cohen, Shay B., Robert J. Simmons, and Noah A. Smith (2011). "Products of Weighted Logic Programs". In: *Theory and Practice of Logic Programming* 11.

 URL: https://link.springer.com/chapter/10.1007/978-3-540-89982-2_18.
- Cohn, Paul Moritz (1969). *Universal Algebra*. Vol. 6. Harper and Row. URL: https://link.springer.com/book/10.1007/978-94-009-8399-1.
- Collins, Michael John (1996). "A New Statistical Parser Based on Bigram Lexical Dependencies". In: *Proceedings of the Conference of the Association for Computational Linguistics*. URL: https://www.aclweb.org/anthology/P96-1025.
- Colmerauer, Alain and Philippe Roussel (1996). "The Birth of Prolog". In: *History of Programming languages*. Vol. II. Association for Computing Machinery.

 URL: https://doi.org/10.1145/234286.1057820.

- Colón, Michael and Henny Sipma (2001). "Synthesis of Linear Ranking Functions". In: *Proceedings of the International Conference on Tools and Algorithms for the Construction and Analysis of Systems*. Vol. 2031. Springer. URL: https://doi.org/10.1007/3-540-45319-9_6.
- Cousot, Patrick and Radhia Cousot (1977). "Abstract Interpretation: A Unified Lattice Model for Static Analysis of Programs by Construction or Approximation of Fixpoints". In: Conference Record of the ACM Symposium on Principles of Programming Languages. URL: https://doi.org/10.1145/512950.512973.
- Dawson, Steven, C. R. Ramakrishnan, and David S. Warren (1996). "Practical Program Analysis Using General Purpose Logic Programming Systems—a Case Study". In: *Proceedings of the ACM SIGPLAN Conference on Programming Language Design and Implementation*. URL: https://doi.org/10.1145/231379.231399.
- Deshpande, Amol, Zachary G. Ives, and Vijayshankar Raman (2007). "Adaptive Query Processing". In: Foundations and Trends® in Databases 1.1. URL: https://doi.org/10.1561/1900000001.
- Dial, Robert B. (1969). "Algorithm 360: Shortest-Path Forest with Topological Ordering [H]". In: Communications of the Association for Computing Machinery 12.11. URL: https://doi.org/10.1145/363269.363610.
- Dijkstra, Edsger W. (1959). "A Note on Two Problems in Connexion with Graphs".

 In: Numerische Mathematik 1. URL: https://doi.org/10.1007/BF01386390.
- Dolan, Stephen (2013). "Fun with Semirings: A Functional Pearl on the Abuse of Linear Algebra". In: *Proceedings of the ACM SIGPLAN International Con-*

- ference on Functional Programming. Association for Computing Machinery.

 URL: https://doi.org/10.1145/2500365.2500613.
- Dozat, Timothy and Christopher D. Manning (2017). "Deep Biaffine Attention for Neural Dependency Parsing". In: *Proceedings of the International Conference on Learning Representations*. URL: https://openreview.net/forum?id=Hk95PK9le.
- Durrett, Greg and Dan Klein (2015). "Neural CRF Parsing". In: Proceedings of the Annual Meeting of the Association for Computational Linguistics and the International Joint Conference on Natural Language Processing. URL: https://aclanthology.org/P15-1030.
- Dyer, Chris (2013). "Minimum Error Rate Training and the Convex Hull Semiring". In: Computing Research Repository abs/1307.3675. arXiv: 1307.3675. URL: http://arxiv.org/abs/1307.3675.
- Eisner, Jason (2001). "Expectation Semirings: Flexible EM for Finite-State Transducers". In: *Proceedings of the ESSLLI Workshop on Finite-State Methods in Natural Language Processing*. Helsinki. URL: http://cs.jhu.edu/~jason/papers/#eisner-2001-fsmnlp.
- Eisner, Jason (2002). "Parameter Estimation for Probabilistic Finite-State Transducers". In: *Proceedings of the Conference of the Association for Computational Linguistics*. URL: https://www.aclweb.org/anthology/P02-1001.
- Eisner, Jason (2023). "Time-and-Space-Efficient Weighted Deduction". In: *Transactions of the Association for Computational Linguistics*. URL: http://cs.jhu.edu/~jason/papers/#eisner-2023-tacl.

- Eisner, Jason and John Blatz (2007). "Program Transformations for Optimization of Parsing Algorithms and Other Weighted Logic Programs". In: *Proceedings of the Conference on Formal Grammar*. URL: http://cs.jhu.edu/~jason/papers/#eisner-blatz-2007.
- Eisner, Jason and Nathaniel W. Filardo (2011). "Dyna: Extending Datalog For Modern AI". In: *Datalog Reloaded*. Vol. 6702. Lecture Notes in Computer Science. Springer. URL: https://link.springer.com/chapter/10.1007/978-3-642-24206-9_11.
- Eisner, Jason, Eric Goldlust, and Noah A. Smith (2004). "Dyna: A Declarative Language for Implementing Dynamic Programs". In: *Proceedings of the Conference of the Association for Computational Linguistics*. URL: https://www.aclweb.org/anthology/H05-1036/.
- Eisner, Jason, Eric Goldlust, and Noah A. Smith (2005). "Compiling Comp Ling: Weighted Dynamic Programming and the Dyna Language". In: Proceedings of the Joint Conference on Human Language Technology Conference and Empirical Methods in Natural Language Processing. URL: https://aclanthology.org/H05-1036/.
- Eisner, Jason, Michael Kornbluh, Gordon Woodhull, Raymond Buse, Samuel Huang, Constantinos Michael, and George Shafer (2006). "Visual Navigation Through Large Directed Graphs and Hypergraphs". In: *Proceedings of the IEEE Symposium on Information Visualization*. URL: http://cs.jhu.edu/~jason/papers/#DYNASTY-2006.

- Eisner, Jason and Giorgio Satta (1999a). "Efficient Parsing for Bilexical Context-Free Grammars and Head Automaton Grammars". In: *Proceedings of the* Conference of the Association for Computational Linguistics. URL: https://www.aclweb.org/anthology/P99-1059.
- Eisner, Jason and Giorgio Satta (1999b). "Efficient Parsing for Bilexical Context-Free Grammars and Head-Automaton Grammars". In: *Proceedings of the* Conference of the Association for Computational Linguistics. URL: https://aclanthology.org/P99-1059/.
- Eisner, Jason and Giorgio Satta (2000). "A Faster Parsing Algorithm for Lexicalized Tree-Adjoining Grammars". In: *Proceedings of the International Workshop on Tree Adjoining Grammar and Related Frameworks*. URL: https://aclanthology.org/W00-2011/.
- Eisner, Jason M. (1996). "Three New Probabilistic Models for Dependency Parsing: An Exploration". In: *Proceedings of the Conference of the Association for Computational Linguistics*. URL: https://www.aclweb.org/anthology/C96-1058.
- Elidan, Gal, Ian McGraw, and Daphne Koller (2006). "Residual Belief Propagation: Informed Scheduling for Asynchronous Message Passing". In: *Proceedings of the Conference on Uncertainty in Artificial Intelligence*. URL: https://arxiv.org/abs/1206.6837.
- Esparza, Javier, Stefan Kiefer, and Michael Luttenberger (2007). "An Extension of Newton's Method to ω-Continuous Semirings". In: *Proceedings of the International Conference on Developments in Language Theory*. URL: https://doi.org/10.1007/978-3-540-73208-2_17.

- Esparza, Javier, Stefan Kiefer, and Michael Luttenberger (2010). "Newtonian Program Analysis". In: *Journal of the ACM* 57.6. URL: https://doi.org/10.1145/1857914.1857917.
- Etessami, Kousha and Mihalis Yannakakis (2009). "Recursive Markov Chains, Stochastic Grammars, and Monotone Systems of Nonlinear Equations". In:

 Journal of the Association for Computing Machinery 56.1. URL: https://doi.org/10.1145/1462153.1462154.
- Felzenszwalb, P. F. and D. McAllester (2007). "The Generalized A* Architecture".

 In: Journal of Artificial Intelligence Research 29. URL: http://www.jair.org/papers/paper2187.html.
- Filardo, Nathaniel Wesley (2017). "Dyna 2: Towards a General Weighted Logic Language". PhD thesis. Johns Hopkins University. URL: https://jscholarship.library.jhu.edu/bitstream/handle/1774.2/58686/FILARDO-DISSERTATION-2017.pdf.
- Filardo, Nathaniel Wesley and Jason Eisner (2012). "A Flexible Solver for Finite Arithmetic Circuits". In: *International Conference on Logic Programming*. Vol. 17. Leibniz International Proceedings in Informatics. URL: http://cs.jhu.edu/~jason/papers/#filardo-eisner-2012-iclp.
- Floyd, Robert W. (1962). "Algorithm 97: Shortest Path". In: Communications of the Association for Computing Machinery 5.6. URL: https://doi.org/10.1145/367766.368168.
- Ford Jr., Lester Randolph (1956). Network Flow Theory. URL: https://link.springer.com/chapter/10.1007/0-387-22633-8_9.

- Francis-Landau, Matthew, Tim Vieira, and Jason Eisner (2020). "Evaluation of Logic Programs with Built-Ins and Aggregation: A Calculus for Bag Relations". In: *International Workshop on Rewriting Logic and Its Applications*.

 URL: http://cs.jhu.edu/~jason/papers/#francislandau-vieira-eisner-2020-wrla.
- Frigo, Matteo and Steven G. Johnson (1998). "FFTW: An Adaptive Software Architecture for the FFT". In: *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing*. Vol. 3. URL: https://doi.org/10.1109%2FICASSP.1998.681704.
- Frühwirth, Thom (1998). "Theory and Practice of Constraint Handling Rules".

 In: *The Journal of Logic Programming* 37.1. URL: http://www.sciencedirect.

 com/science/article/pii/S0743106698100055.
- Galler, Bernard A. and Michael J. Fischer (1964). "An Improved Equivalence Algorithm". In: Communications of the Association for Computing Machinery 7.5. URL: https://doi.org/10.1145/364099.364331.
- Gallo, Giorgio, Giustino Longo, Stefano Pallottino, and Sang Nguyen (1993). "Directed Hypergraphs and Applications". In: *Discrete Applied Mathematics* 42.2.

 URL: https://www.sciencedirect.com/science/article/pii/0166218X9390045P.
- Ganzinger, Harald and David McAllester (2002). "Logical Algorithms". In: Logic

 Programming. URL: https://link.springer.com/chapter/10.1007/3-540-45619-8_15.
- Gergatsoulis, Manolis and Maria Katzouraki (1994). "Unfold/Fold Transformations For Definite Clause Programs". In: *Proceedings of Programming*

- Language Implementation and Logic Programming. URL: https://doi.org/10.1007/3-540-58402-1_24.
- Gildea, Daniel (2011). "Grammar Factorization by Tree Decomposition". In:

 *Computational Linguistics 37.1. URL: https://www.aclweb.org/anthology/J11
 1008.
- Gildea, Daniel (2020). "Efficient Outside Computation". In: Computational Linguistics 46.4. URL: https://aclanthology.org/2020.cl-4.2.
- Gilroy, Sorcha, Adam Lopez, and Sebastian Maneth (2017). "Parsing Graphs with Regular Graph Grammars". In: *Proceedings of the Joint Conference on Lexical and Computational Semantics*. URL: https://www.aclweb.org/anthology/S17-1024.
- Gómez-Rodríguez, Carlos (2010). Parsing Schemata for Practical Text Analysis.

 Vol. 1. Mathematics, Computing, Language, and Life: Frontiers in Mathematical Linguistics and Language Theory. Imperial College Press. URL: http://www.worldscientific.com/worldscibooks/10.1142/p714.
- Gómez-Rodríguez, Carlos, Jesús Vilares, and Miguel A. Alonso (2009). "A Compiler for Parsing Schemata". In: *Software—Practice & Experience* 39.5. URL: https://doi.org/10.1002/spe.904.
- Gondran, Michel and Michel Minoux (2008). *Graphs, Dioids and Semirings: New Models and Algorithms*. Vol. 41. Springer Science & Business Media. URL: https://link.springer.com/book/10.1007/978-0-387-75450-5.
- Goodman, Joshua (1999). "Semiring Parsing". In: Computational Linguistics 25.4. URL: https://www.aclweb.org/anthology/J99-4004.

- Goodman, Joshua T. (1998). "Parsing Inside-out". PhD thesis. Harvard University. URL: https://arxiv.org/abs/cmp-lg/9805007.
- Green, Todd J., Grigoris Karvounarakis, and Val Tannen (2007). "Provenance Semirings". In: *Proceedings of the ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems*. URL: https://doi.org/10.1145/1265530.1265535.
- Green, Todd J. and Val Tannen (2017). "The Semiring Framework for Database Provenance". In: *Proceedings of the ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems*. URL: https://doi.org/10.1145/3034786. 3056125.
- Griewank, Andreas and Andrea Walther (2008). Evaluating Derivatives: Principles and Techniques of Algorithmic Differentiation. 2nd ed. SIAM. URL: https://doi.org/10.1137/1.9780898717761.
- Hall, Keith B. (2005). "Best-first Word-lattice Parsing: Techniques for Integrated Syntactic Language Modeling". PhD thesis. Brown University. URL: https://www.proquest.com/openview/eeb76c0001e6d80db7c72be0ea977440/.
- Harris, Charles R., K. Jarrod Millman, St'efan J. van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fern'andez del R'io, Mark Wiebe, Pearu Peterson, Pierre G'erard-Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and

- Travis E. Oliphant (2020). "Array Programming with NumPy". In: *Nature* 585.7825. URL: https://doi.org/10.1038/s41586-020-2649-2.
- Herbrand, Jacques (1930). "Recherches sur la Théorie de la Demonstration".

 PhD thesis. Université de Paris. URL: http://www.numdam.org/issue/THESE_1930_
 __110__1_0.pdf.
- Hong, Juneki and Liang Huang (2018). "Linear-time Constituency Parsing with RNNs and Dynamic Programming". In: *Proceedings of the Annual Meeting of the Association for Computational Linguistics*. URL: https://aclanthology.org/P18-2076.
- Hopcroft, John E., Rajeev Motwani, and Jeffrey D. Ullman (2007). *Introduction to Automata Theory, Languages, and Computation*. 3rd Edition. Addison—Wesley. URL: https://dl.acm.org/doi/book/10.5555/1196416.
- Huang, Liang (2008). "Advanced Dynamic Programming in Semiring and Hypergraph Frameworks". In: COLING: Advanced Dynamic Programming in Computational Linguistics: Theory, Algorithms and Applications Tutorial notes. URL: https://aclanthology.org/C08-5001.
- Huang, Liang and David Chiang (2005). "Better k-best Parsing". In: Proceedings of the International Workshop on Parsing Technology. URL: https://aclanthology.org/W05-1506.
- Huang, Liang and Kenji Sagae (2010). "Dynamic Programming for Linear-Time Incremental Parsing". In: Proceedings of the Conference of the Association for Computational Linguistics. URL: https://aclanthology.org/P10-1110.

- Huang, Liang, Hao Zhang, and Daniel Gildea (2005). "Machine Translation as Lexicalized Parsing with Hooks". In: *International Workshop on Parsing Technologies*. URL: https://www.aclweb.org/anthology/W05-1507.
- Hutter, Frank (2009). "Automated Configuration of Algorithms for Solving Hard Computational Problems". PhD thesis. University Of British Columbia (Vancouver). URL: https://www.cs.ubc.ca/~hutter/papers/Hutter09PhD.pdf.
- Iwano, Kazuo and Kenneth Steiglitz (1990). "A Semiring on Convex Polygons and Zero-Sum Cycle Problems". In: *SIAM Journal on Computing* 19.5. URL: https://doi.org/10.1137/0219061.
- Jelinek, Frederick (1985). "Markov Source Modeling of Text Generation". In:

 NATO Advanced Study Institute: Impact of Processing Techniques on Communication. URL: https://link.springer.com/chapter/10.1007/978-94-009-5113-6_28.
- Jiang, Jiarong, Adam Teichert, Hal Daumé III, and Jason Eisner (2012). "Learned Prioritization for Trading Off Accuracy and Speed". In: *Advances in Neural Information Processing Systems*. URL: http://cs.jhu.edu/~jason/papers/#jianget-al-2012-nips.
- Johnson, Mark (1989). "The Computational Complexity of Tomita's Algorithm".

 In: International Workshop on Parsing Technologies. URL: https://www.aclweb.org/anthology/W89-0221.
- Johnson, Mark (1998). "Finite-state Approximation of Constraint-based Grammars using Left-corner Grammar Transforms". In: Annual Meeting of the

- Association for Computational Linguistics and International Conference on Computational Linguistics. URL: https://aclanthology.org/P98-1101.
- Johnson, Mark (2007). "Transforming Projective Bilexical Dependency Grammars into efficiently-parsable CFGs with Unfold-Fold". In: *Proceedings of the Conference of the Association for Computational Linguistics*. URL: https://www.aclweb.org/anthology/P07-1022.
- Johnson, Mark and Brian Roark (2000). "Compact Non-left-recursive Grammars using the Selective Left-corner Transform and Factoring". In: *Proceedings of The International Conference on Computational Linguistics*. URL: https://aclanthology.org/C00-1052.
- Jurafsky, Dan and James H. Martin (2020). Speech and Language Processing:

 An Introduction to Natural Language Processing, Computational Linguistics,
 and Speech Recognition, Third Edition. Third. (Draft). URL: https://web.
 stanford.edu/~jurafsky/slp3/.
- Kasami, Tadao (1965). An Efficient Recognition and Syntax Algorithm for Context-free Languages. Tech. rep. AF-CRL-65-758. Air Force Cambridge Research Laboratory. URL: https://www.ideals.illinois.edu/handle/2142/74304.
- Kawamura, Tadashi and Tadashi Kanamori (1990). "Preservation of Stronger Equivalence in Unfold/Fold Logic Program Transformation". In: *Theoretical Computer Science* 75.1. URL: https://www.sciencedirect.com/science/article/pii/030439759090065P.

- Kim, Yoon, Carl Denton, Luong Hoang, and Alexander M. Rush (2017). "Structured Attention Networks". In: *Proceedings of the International Conference on Learning Representations*. URL: https://openreview.net/forum?id=HkE0Nvqlg.
- Kleene, Stephen C. (1956). "Representation of Events in Nerve Nets and Finite Automata". In: *Automata Studies* 34. URL: https://apps.dtic.mil/sti/pdfs/ADA596138.pdf.
- Klein, Dan and Chris D. Manning (2001). "Parsing and Hypergraphs". In: *Proceedings of the Seventh International Workshop on Parsing Technologies*. URL: https://aclanthology.org/W01-1812.
- Knight, Kevin (1989). "Unification: A Multidisciplinary Survey". In: ACM Computing Surveys 21.1. URL: https://doi.org/10.1145/62029.62030.
- Knuth, Donald E. (1977). "A generalization of Dijkstra's algorithm". In: *Information Processing Letters* 6.1. URL: https://www.sciencedirect.com/science/article/pii/0020019077900023.
- Koo, Terry and Michael Collins (2010). "Efficient Third-Order Dependency Parsers". In: *Proceedings of the Conference of the Association for Computational Linguistics*. URL: https://www.aclweb.org/anthology/P10-1001.
- Kowalski, Robert (1979). "Algorithm = Logic + Control". In: Communications of the Association for Computing Machinery 22.7. URL: https://dl.acm.org/citation.cfm?id=359136.
- Kuhlmann, Marco (2013). "Mildly Non-Projective Dependency Grammar". In: Computational Linguistics 39.2. URL: https://www.aclweb.org/anthology/J13-2004.

- Kuhlmann, Marco, Carlos Gómez-Rodríguez, and Giorgio Satta (2011). "Dynamic Programming Algorithms for Transition-Based Dependency Parsers". In: Proceedings of the Annual Meeting of the Association for Computational Linguistics: Human Language Technologies. URL: https://aclanthology.org/P11-1068.
- Kuhlmann, Marco and Giorgio Satta (2014). "A New Parsing Algorithm for Combinatory Categorial Grammar". In: *Transactions of the Association for Computational Linguistics* 2. URL: https://aclanthology.org/Q14-1032.
- Kuhlmann, Marco, Giorgio Satta, and Peter Jonsson (2018). "On the Complexity of CCG Parsing". In: *Computational Linguistics* 44.3. URL: https://www.aclweb.org/anthology/J18-3004.
- Kuich, Werner (1997). "Semirings and Formal Power Series: Their Relevance to Formal Languages and Automata". In: *Handbook of Formal Languages: Word, Language, Grammar*. Vol. 1. Springer. URL: https://doi.org/10.1007/978-3-642-59136-5_9.
- Lafferty, John, Andrew McCallum, and Fernando Pereira (2001). "Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data". In: *Proceedings of the International Conference on Machine Learning*.

 URL: https://dl.acm.org/doi/10.5555/645530.655813.
- Lange, Martin and Hans Leiß (2009). "To CNF or not to CNF? An Efficient Yet Presentable Version of the CYK Algorithm". In: *Informatica Didactica* 8. URL: http://www.informatica-didactica.de/cmsmadesimple/index.php?page=LangeLeiss2009.

- Lauritzen, S. L. and D. J. Spiegelhalter (1988). "Local Computations with Probabilities on Graphical Structures and their Application to Expert Systems". In:

 Journal of the Royal Statistical Society. Series B (Methodological) 50.2. URL:

 https://rss.onlinelibrary.wiley.com/doi/10.1111/j.2517-6161.1988.tb01721.x.
- Lee, Kenton, Mike Lewis, and Luke Zettlemoyer (2016). "Global Neural CCG Parsing with Optimality Guarantees". In: *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. URL: https://aclanthology.org/D16-1262.
- Lehmann, Daniel J. (1977). "Algebraic Structures for Transitive Closure". In:

 Theoretical Computer Science 4.1. URL: https://www.sciencedirect.com/science/article/pii/0304397577900561.
- Leis, Viktor, Bernharde Radke, Andrey Gubichev, Alfons Kemper, and Thomas Neumann (2017). "Cardinality Estimation Done Right: Index-Based Join Sampling". In: *Proceedings of the Conference on Innovative Data Systems Research*. URL: https://www.cidrdb.org/cidr2017/papers/p9-leis-cidr17.pdf.
- Levenshtein, Vladimir I. (1966). "Binary Codes Capable of Correcting Deletions, Insertions, and Reversals". In: *Soviet Physics Doklady* 8 (10). URL: https://nymity.ch/sybilhunting/pdf/Levenshtein1966a.pdf.
- Li, Zhifei and Jason Eisner (2009). "First- and Second-Order Expectation Semirings with Applications to Minimum-Risk Training on Translation Forests".

 In: Proceedings of the Conference on Empirical Methods in Natural Language

 Processing. URL: https://www.aclweb.org/anthology/D09-1005.

- Lopez, Adam (2009). "Translation as Weighted Deduction". In: *Proceedings of the Conference of the European Association for Computational Linguistics*.

 URL: https://www.aclweb.org/anthology/E09-1061.
- Łukasiewicz, Jan (1920). "O Logice Trojwartosciowej". In: Ruch Filozoficzny 5.
- Luttenberger, Michael (2009). "Solving Polynomial Systems on Semirings: A Generalization of Newton's Method". PhD thesis. TU München. URL: https://d-nb.info/100090086X/34.
- Mächler, Martin (2012). Accurately Computing $\log(1-\exp(-|a|))$ Assessed by the Rmpfr Package. Tech. rep. URL: https://cran.r-project.org/web/packages/Rmpfr/vignettes/log1mexp-note.pdf.
- Maher, Michael J. (1987). Correctness of a Logic Program Transformation System.

 Tech. rep. RC13496. IBM Research Report, T.J. Watson Research. URL: https://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.92.2964.
- Mannila, Heikki and Esko Ukkonen (1986). "On the Complexity of Unification Sequences". In: *Proceedings of the International Conference on Logic Programming*. URL: https://link.springer.com/chapter/10.1007/3-540-16492-8_69.
- Marcus, Mitchell P., Beatrice Santorini, and Mary Ann Marcinkiewicz (1993).

 "Building a Large Annotated Corpus of English: The Penn Treebank". In:

 Computational Linguistics 19.2. URL: https://www.aclweb.org/anthology/J93-2004.
- Martelli, Alberto and Ugo Montanari (1982). "An Efficient Unification Algorithm". In: *ACM Transactions on Programming Languages and Systems* 4.2. URL: https://dl.acm.org/doi/10.1145/357162.357169.

- McAllester, David (2002). "On the Complexity Analysis of Static Analyses". In: Journal of the ACM 49.4. URL: https://doi.org/10.1145/581771.581774.
- Melamed, I. Dan (2003). "Multitext Grammars and Synchronous Parsers". In:

 Proceedings of the Human Language Technology Conference of the North

 American Chapter of the Association for Computational Linguistics. URL:

 https://www.aclweb.org/anthology/N03-1021.
- Melamed, I. Dan, Giorgio Satta, and Benjamin Wellington (2004). "Generalized Multitext Grammars". In: *Proceedings of the Annual Meeting of the Association for Computational Linguistics*. Barcelona, Spain. URL: https://aclanthology.org/P04-1084.
- Meurer, Aaron, Christopher P. Smith, Mateusz Paprocki, Ondřej Čertík, Sergey B. Kirpichev, Matthew Rocklin, AMiT Kumar, Sergiu Ivanov, Jason K. Moore, Sartaj Singh, Thilina Rathnayake, Sean Vig, Brian E. Granger, Richard P. Muller, Francesco Bonazzi, Harsh Gupta, Shivam Vats, Fredrik Johansson, Fabian Pedregosa, Matthew J. Curry, Andy R. Terrel, Štěpán Roučka, Ashutosh Saboo, Isuru Fernando, Sumith Kulal, Robert Cimrman, and Anthony Scopatz (2017). "SymPy: Symbolic Computing in Python". In: *PeerJ Computer Science* 3. URL: https://doi.org/10.7717/peerj-cs.103.
- Minnen, Guido (1996). "Magic for Filter Optimization in Dynamic Bottom-Up Processing". In: *Proceedings of the Annual Meeting of the Association for Computational Linguistics*. URL: https://aclanthology.org/P96-1033.

- Mohri, Mehryar (1997). "Finite-State Transducers in Language and Speech Processing". In: Computational Linguistics 23.2. URL: https://www.aclweb.org/anthology/J97-2003.
- Mohri, Mehryar (2002). "Semiring Frameworks and Algorithms for Shortest-Distance Problems". In: *Journal of Automata, Languages and Combinatorics* 7.3. URL: https://cs.nyu.edu/~mohri/pub/jalc.pdf.
- Mohri, Mehryar (2009). "Weighted Automata Algorithms". In: *Handbook of Weighted Automata*. Springer Berlin Heidelberg. URL: https://doi.org/10.1007/978-3-642-01492-5_6.
- Moore, Edward F. (1959). "The Shortest Path through a Maze". In: *Proceedings* of an International Symposium on the Theory of Switching. URL: https://books.google.com/books/about/The_Shortest_Path_Through_a_Maze.html?id=IVZBHAAACAAJ.
- Moura, Leonardo Mendonça de and Nikolaj Bjørner (2008). "Z3: An Efficient SMT Solver". In: *Proceedings of the International Conference on Tools and Algorithms for the Construction and Analysis of Systems*. URL: https://doi.org/10.1007/978-3-540-78800-3_24.
- Nederhof, Mark-Jan and Ricardo Sánchez-Sáez (2011). "Parsing of Partially Bracketed Structures for Parse Selection". In: *International Workshop on Parsing Technologies*. URL: https://www.aclweb.org/anthology/W11-2928.
- Nederhof, Mark-Jan and Giorgio Satta (2008). "Computing Partition Functions of PCFGs". In: *Research on Language and Computation*. URL: https://link.springer.com/content/pdf/10.1007/s11168-008-9052-8.pdf.

- Nederhof, Mark-Jan and Giorgio Satta (2011). "Prefix Probability for Probabilistic Synchronous Context-Free Grammars". In: *Proceedings of the Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*. URL: https://www.aclweb.org/anthology/P11-1047.
- Nivre, Joakim (2003). "An Efficient Algorithm for Projective Dependency Parsing". In: *Proceedings of the International Conference on Parsing Technologies*.

 URL: https://aclanthology.org/w03-3017.
- Norvig, Peter (1991). "Correcting a Widespread Error in Unification Algorithms".

 In: Software: Practice and Experience 21.2. URL: http://norvig.com/unifybug.pdf.
- Och, Franz Josef (2003). "Minimum Error Rate Training in Statistical Machine Translation". In: *Proceedings of the Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics. URL: http://www.aclweb.org/anthology/P03-1021.
- Overton, David (2003). "Precise and Expressive Mode Systems for Typed Logic Programming Languages". PhD thesis. University of Melbourne. URL: https://www.semanticscholar.org/paper/Precise-and-Expressive-Mode-Systems-for-Typed-Logic-Overton/53a9c78610eda8ab617fa47d1f60d4c6c9f31957.
- Paszke, Adam, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer (2017). "Automatic Differentiation in PyTorch". In: *Proceedings of the NeurIPS Autodiff Workshop*. URL: http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf.

- Pereira, Fernando C. N. and David H. D. Warren (1983). "Parsing as Deduction". In: *Proceedings of the Conference of the Association for Computational Linguistics*. URL: https://aclanthology.org/P83-1021.
- Pettorossi, Alberto and Maurizio Proietti (1994). "Transformation of Logic Programs: Foundations and Techniques". In: *The Journal of Logic Programming* 19–20. URL: https://www.sciencedirect.com/science/article/pii/0743106694900280.
- Plotkin, Gordon D (2004). "The Origins of Structural Operational Semantics".

 In: The Journal of Logic and Algebraic Programming 60-61. URL: https://www.sciencedirect.com/science/article/pii/S1567832604000268.
- Plotkin, Gordon D. (1970). "A Note on Inductive Generalization". In: *Machine Intelligence*. URL: https://homepages.inf.ed.ac.uk/gdp/publications/MI5_note_ind_gen.pdf.
- Podelski, Andreas and Andrey Rybalchenko (2004). "A Complete Method for the Synthesis of Linear Ranking Functions". In: *Proceedings of the International Conference on Verification, Model Checking, and Abstract Interpretation*. Vol. 2937. Springer. URL: https://doi.org/10.1007/978-3-540-24622-0_20.
- Prieditis, Armand E. (1993). "Machine Discovery of Effective Admissible Heuristics". In: *Machine Learning* 12.1–3. URL: https://doi.org/10.1007/BF00993063.
- Purdom Jr., Paul Walton (1970). "A Transitive Closure Algorithm". In: *BIT Numerical Mathematics* 10 (1). URL: https://link.springer.com/article/10.1007/BF01940892.

- Qi, Peng, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D. Manning (2020). "Stanza: A Python Natural Language Processing Toolkit for Many Human Languages". In: Proceedings of the Annual Meeting of the Association for Computational Linguistics: System Demonstrations. URL: https://aclanthology.org/2020.acl-demos.14.
- Rabiner, Lawrence R (1989). "A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition". In: *Proceedings of the IEEE* 77.2. URL: https://ieeexplore.ieee.org/document/18626.
- Ramakrishnan, Raghu (1991). "Magic Templates: A Spellbinding Approach To Logic Programs". In: *Journal of Logic Programming* 11.3–4. URL: https://doi.org/10.1016/0743-1066(91)90026-L.
- Rastogi, Pushpendre, Ryan Cotterell, and Jason Eisner (2016). "Weighting Finite-State Transductions With Neural Context". In: *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. URL: https://aclanthology.org/N16-1076.
- Reps, Thomas W. (1993). "Demand Interprocedural Program Analysis Using Logic Databases". In: *Applications of Logic Databases*. The Kluwer International Series in Engineering and Computer Science. URL: https://link.springer.com/chapter/10.1007/978-1-4615-2207-2_8.
- Reps, Thomas W., Emma Turetsky, and Prathmesh Prabhu (2017). "Newtonian Program Analysis via Tensor Product". In: *ACM Transactions on Programming Language Systems* 39.2. URL: https://doi.org/10.1145/3024084.

- Reynolds, John C. (1970). "Transformational Systems and the Algebraic Structure of Atomic Formulas". In: *Machine Intelligence*. URL: http://www.cs.cmu.edu/afs/cs/user/jcr/ftp/transysalg.pdf.
- Robinson, John Alan (1965). "A Machine-Oriented Logic Based on the Resolution Principle". In: *Journal of the ACM* 12.1. URL: http://doi.acm.org/10.1145/321250.321253.
- Rosenkrantz, Daniel J. and Philip M. Lewis (1970). "Deterministic Left Corner Parsing". In: *Proceedings of The Annual Symposium on Switching and Automata Theory*. URL: https://ieeexplore.ieee.org/document/4569645.
- Roy, Bernard (1959). "Transitivité et Connexité". In: Comptes rendus de l'Académie des Sciences 249.2. URL: https://gallica.bnf.fr/ark:/12148/bpt6k3201c/f222.image.
- Roychoudhury, Abhik, K. Narayan Kumar, C. R. Ramakrishnan, and I. V. Ramakrishnan (1999). "A Parameterized Unfold/Fold Transformation Framework for Definite Logic Programs". In: *Principles and Practice of Declarative Programming*. URL: https://www3.cs.stonybrook.edu/~cram/Papers/RKRR_PPDP99/paper.pdf.
- Rush, Alexander (2020). "Torch-Struct: Deep Structured Prediction Library".

 In: Proceedings of the Annual Meeting of the Association for Computational

 Linguistics. URL: https://aclanthology.org/2020.acl-demos.38.
- Russell, Stuart J. and Peter Norvig (2020). Artificial Intelligence: A Modern

 Approach (4th Edition). Pearson. URL: http://aima.cs.berkeley.edu/.

- Sarawagi, Sunita and William W. Cohen (2004). "Semi-Markov Conditional Random Fields for Information Extraction". In: *Advances in Neural Information Processing Systems*. URL: https://proceedings.neurips.cc/paper/2004/hash/eb06b9db06012a7a4179b8f3cb5384d3-Abstract.html.
- Schlund, Maximilian (2016). "Algebraic Systems of Fixpoint Equations over Semirings: Theory and Applications". PhD thesis. Technical University Munich. URL: https://mediatum.ub.tum.de/doc/1273784/1273784.pdf.
- Schlund, Maximilian, Michal Terepeta, and Michael Luttenberger (2013). "Putting Newton into Practice: A Solver for Polynomial Equations over Semirings".

 In: Proceedings of the International Conference on Logic for Programming,

 Artificial Intelligence, and Reasoning. URL: https://doi.org/10.1007/978-3-642-45221-5_48.
- Shafer, Glenn and Prakash P. Shenoy (1990). "Probability Propagation". In:

 Annals of Mathematics and Artificial Intelligence 2 (1). URL: https://doi.org/
 10.1007/BF01531015.
- Shi, Tianze, Liang Huang, and Lillian Lee (2017). "Fast(er) Exact Decoding and Global Training for Transition-Based Dependency Parsing via a Minimal Feature Set". In: *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. URL: https://aclanthology.org/D17-1002.
- Shieber, Stuart M., Yves Schabes, and Fernando C. N. Pereira (1995). "Principles and Implementation of Deductive Parsing". In: *Journal of Logic Programming* 24.1–2. URL: http://www.eecs.harvard.edu/~shieber/Biblio/Papers/infer.pdf.

- Sikkel, Klaus (1997). Parsing Schemata: A Framework for Specification and Analysis of Parsing Algorithms. Texts in Theoretical Computer Science. Springer-Verlag. URL: https://link.springer.com/book/10.1007/978-3-642-60541-3.
- Smith, Daniel G. A. and Johnnie Gray (2018). "opt_einsum A Python Package for Optimizing Contraction Order for Einsum-like Expressions". In: *Journal of Open Source Software* 3.26. URL: https://doi.org/10.21105/joss.00753.
- Stern, Mitchell, Jacob Andreas, and Dan Klein (2017). "A Minimal Span-Based Neural Constituency Parser". In: *Proceedings of the Annual Meeting of the Association for Computational Linguistics*. URL: https://aclanthology.org/P17-1076.
- Stolcke, Andreas (1995). "An Efficient Probabilistic Context-Free Parsing Algorithm that Computes Prefix Probabilities". In: *Computational Linguistics* 21.2. URL: https://www.aclweb.org/anthology/J95-2002.
- Strassen, Volker (1969). "Gaussian Elimination is Not Optimal". In: Numerische Mathematik 13.4. URL: https://doi.org/10.1007/BF02165411.
- Tamaki, Hisao and Taisuke Sato (1984). "Unfold/Fold Transformation of Logic Programs". In: *Proceedings of the International Logic Programming Conference*. URL: https://archive.computerhistory.org/resources/access/text/2022/11/102785625-05-01-acc.pdf.
- Tarjan, Robert (1972). "Depth-First Search and Linear Graph Algorithms". In: SIAM Journal on Computing 1.2. URL: https://doi.org/10.1137/0201010.

- Tarjan, Robert Endre (1981). "A Unified Approach to Path Problems". In: Journal of the Association for Computing Machinery 28.3. URL: https://doi.org/10.1145/322261.322272.
- Teitelbaum, Ray (1973). "Context-free Error Analysis by Evaluation of Algebraic Power Series". In: *Proceedings of the ACM Symposium on Theory of Computing*. URL: https://dl.acm.org/doi/pdf/10.1145/800125.804050.
- Theano Development Team (2016). "Theano: A Python Framework for Fast Computation of Mathematical Expressions". In: Computing Research Repository abs/1605.02688. URL: http://arxiv.org/abs/1605.02688.
- Tomita, Masaru (1985). "An Efficient Context-Free Parsing Algorithm for Natural Languages". In: *Proceedings of the International Joint Conference on Artificial Intelligence*. URL: https://dl.acm.org/doi/10.5555/1623611.1623625.
- Tomita, Masaru (1986). "An Efficient Word Lattice Parsing Algorithm for Continuous Speech Recognition". In: *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing*. URL: https://ieeexplore.ieee.org/abstract/document/1168663.
- Tomita, Masaru (1991). *Generalized LR Parsing*. Springer. URL: https://www.springer.com/gp/book/9780792392019.
- Ullman, Jeffrey D. (1988). *Principles of Database and Knowledge-Base Systems*.

 Vol. 1. Computer Science Press. URL: https://www.worldcat.org/oclc/310956623.
- Vieira, Tim, Ryan Cotterell, and Jason Eisner (2021). "Searching for More Efficient Dynamic Programs". In: Findings of the Association for Computational Linguistics. URL: https://arxiv.org/abs/2109.06966.

- Vieira, Tim, Matthew Francis-Landau, Nathaniel Wesley Filardo, Farzad Khorasani, and Jason Eisner (2017). "Dyna: Toward a Self-Optimizing Declarative Language for Machine Learning Applications". In: *Proceedings of the ACM SIGPLAN Workshop on Machine Learning and Programming Languages*).

 URL: http://cs.jhu.edu/~jason/papers/#vieira-et-al-2017.
- Vijay-Shankar, K. and Aravind K. Joshi (1985). "Some Computational Properties of Tree Adjoining Grammars". In: *Proceedings of the Annual Meeting of the Association for Computational Linguistics*. URL: https://aclanthology.org/P85-1011.
- Vijay-Shanker, K. and David J. Weir (1989). "Recognition of Combinatory Categorial Grammars and Linear Indexed Grammars". In: *International Workshop on Parsing Technologies*. URL: https://www.aclweb.org/anthology/W89-0218.
- Vijay-Shanker, K. and David J. Weir (1990). "Polynomial Time Parsing of Combinatory Categorial Grammars". In: *Proceedings of the Conference of the Association for Computational Linguistics*. URL: https://www.aclweb.org/anthology/P90-1001.
- Vijay-Shanker, K. and David J. Weir (1993). "Parsing Some Constrained Grammar Formalisms". In: *Computational Linguistics* 19.4. URL: https://www.aclweb.org/anthology/J93-4002/.
- Viterbi, Andrew J. (1967). "Error Bounds for Convolutional Codes and an Asymptotically Optimum Decoding Algorithm". In: *IEEE Transactions on Information Theory* 13.2. URL: https://doi.org/10.1109/TIT.1967.1054010.

- Warren, David H. D. (1983). An Abstract Prolog Instruction Set. Tech. rep. SRI

 International Artificial Intelligence Center. URL: http://www.ai.sri.com/pubs/files/641.pdf.
- Warshall, Stephen (1962). "A Theorem on Boolean Matrices". In: Journal of the Association for Computing Machinery 9.1. URL: https://doi.org/10.1145/321105.321107.
- Whaley, John, Dzintars Avots, Michael Carbin, and Monica S. Lam (2005).
 "Using Datalog with Binary Decision Diagrams for Program Analysis". In:
 Programming Languages and Systems. URL: https://link.springer.com/
- Whaley, R. Clint, Antoine Petitet, and Jack J. Dongarra (2001). "Automated Empirical Optimizations of Software and the ATLAS Project". In: *Parallel Computing* 27.1. URL: https://www.sciencedirect.com/science/article/abs/pii/S0167819100000879.
- Wingate, David and Kevin D. Seppi (2005). "Prioritization Methods for Accelerating MDP Solvers". In: *Journal of Machine Learning. Research* 6. URL: http://jmlr.org/papers/v6/wingate05a.html.
- Wolfram, Stephen (2003). *The Mathematica Book*. 5th Edition. Wolfram-Media.

 URL: https://www.wolfram.com/books/profile.cgi?id=4939.
- Wu, Dekai (1996). "A Polynomial-Time Algorithm for Statistical Machine Translation". In: *Proceedings of the Conference of the Association for Computational Linguistics*. URL: https://www.aclweb.org/anthology/P96-1021.

- Wu, Shijie and Ryan Cotterell (2019). "Exact Hard Monotonic Attention for Character-Level Transduction". In: *Proceedings of the Annual Meeting of the Association for Computational Linguistics*. URL: https://aclanthology.org/P19-1148.
- Wu, Shijie, Pamela Shapiro, and Ryan Cotterell (2018). "Hard Non-Monotonic Attention for Character-Level Transduction". In: *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. URL: https://aclanthology.org/D18-1473.
- Younger, D. H. (1967). "Recognition and Parsing of Context-Free Languages in Time n^3 ". In: *Information and Control* 10.2. URL: https://core.ac.uk/download/pdf/82305262.pdf.
- Zadeh, L.A. (1965). "Fuzzy Sets". In: *Information and Control* 8.3. URL: https://www.sciencedirect.com/science/article/pii/S001999586590241X.
- Zilberstein, Shlomo (1996). "Using Anytime Algorithms in Intelligent Systems".

 In: AI Magazine 17.3. URL: https://doi.org/10.1609/aimag.v17i3.1232.

Index

not-matches _z (Y), 192	strict partial order, 33
	symmetric, 32
agenda, 81, 133	total order, 33
algebraic path problems, 94 auxiliary definitions, 164	transitive, 32
	well-founded ordering, 33
Banach fixpoint theorem, 69	built-in relation, 14, 25
binary operation, 38	
binary relation, 31	Cauchy sequence, 68
antisymmetric, 32	chart, 73, 133
asymmetric, 32	complete metric space, 68
equivalence, 32	conditional cardinalities, 131
equivalence class, 32	constant folding, 16
infinitely descending chains, 33	Constraint propagation rules, 111
partial order, 33	contraction mapping, 69
preorder, 33	contraction rate, 69
reflexive, 32	cursor trick, 77

depth limit, 115	generalized reversibility test,
derivation, 59	167
height, 60	reversibility test, 163
set of derivations, 59	function, 34
spine, 189	bijection, 34
subtree ordering, 244	injective, 34
value, 59	inverse, 34
yield, 58	surjective, 34
dimension, 88	graph
measure, 176	condensation, 35
set of derivations, 13	strongly connected, 35
derived type, 102	strongly connected components,
distance metric, 68	35
driver, 82	topological sort, 35
driver item, 134	toposort, 35
Dyna 2, 270	graph search problem, 207
fixpoint equation, 60	group, 38
fixpoint iteration, 63	input and output declarations, 143
fold transformation	Input data, 14
basic fold transformation, 162	input type, 102
generalized fold transformation,	item, 10, 12
167	input, 143

intermediate, 143	program, 12
output, 143	fixpoint equations, 60
useful, 120	set of derivations, 59
useless, 120	internal representation, 18
value, 10, 12	semantically equivalent, 152
least upper bound, 68	strongly equivalent, 152
left-corner transformation, 194	support, 102
	program transformation
minimum derivation measure, 183	partial correctness, 166
monoid, 38	safe, 152
commutative, 38	semantics-preserving, 152
idempotent, 38	
multiplicity correction, 31	query, 15
parameter, 104	explicit result stream, 22
parametric expression, 132	implicit result stream, 23
	result, 15
Pareto efficient, 52	
Parikh image, 52	rule, 12
partition, 32	axiom, 14
passenger, 82	body, 12
passengers, 134	grounding, 13
posterior bound, 70	head, 12
prior bound, 69	set of all groundings, 14
prioritization policy, 81	side condition, 14

subgoal, 13	boolean, 44
rule measure, 176	convex sets in a vector space, 53
	counting, 45
semiring, 13, 38	first-order expectation, 46
n-fold product, 39	formal-language, 51
n-fold sum, 39	fuzzy sets, 45
natural order, 66	integer, 45
property	lineage, 50
ω -continuous, 68	log nonnegative real, 46
complete, 67	matrices, 49
naturally ordered, 66	max-plus, 47
property	max-times, 47
K-closed, 40	maximum capacity, 48
closed, 39	min-plus, 47
commutative, 39	min-times, 47
idempotent, 39	Nonnegative real, 46
monotonic, 41	pareto frontier, 52
star operator, 39	provenance, 50
semiring instances	real, 45
K-best, 48	second-order expectation, 46
K-collapsed natural, 46	semilinear sets, 52
Łukasiewicz, 45	sets, 50
arctic, 47	three-valued logic, 45

tropical, 47	to-do list, 213
Viterbi, 47	trace, 189
weighted formal language, 51	transformation sequence, 178
simple, 105	reference program, 178
simple program, 72	truncation, 115
simple type, 107	tupling technique, 157
size parameter, 132	type program, 103
slashed item, 189	unfold transformation
specialized program, 15	basic unfold transformation, 162
speculation transformation, 191	generalized unfold
step operator, 63	transformation, 166
substitution, 19	unification, 19
substitution mapping, 19	update, 81
supremum, 68	valuation, 58
term	valuation function, 12
ground terms, 12	value-preserving bijection, 152
Herbrand universe, 12	variable-elimination technique, 157
nonground term, 18	variables, 13