Syntactic and Semantic Control of Large Language Models via Sequential Monte Carlo

João Loula^{*1} Benjamin LeBrun^{*5} Li Du^{*6} Ben Lipkin¹ Clemente Pasti² Gabriel Grand¹ Tianyu Liu² Yahya Emara² Marjorie Freedman⁸ Jason Eisner⁶ Ryan Cotterell² Vikash Mansinghka^{‡1} Alexander K. Lew^{‡1,7} Tim Vieira^{‡2} Timothy J. O'Donnell^{‡3,4,5} ¹MIT ²ETH Zürich ³McGill ⁴Canada CIFAR AI Chair ⁵Mila ⁶Johns Hopkins ⁷Yale ⁸ISI *co-first authorship, [‡]co-senior authorship

- Code:
 - Static analysis

LSP



```
def signal(context):
    server = LSPDiagnosticServer("python")
    diagnostics = server.get_diagnostics(code)
    errors = [d for d in diagnostics
        if d["severity"] == 1]
        if errors:
            return 0
            return 1
```

- Code:
 - Static analysis
 - Run tests check for errors

```
SciPy
                                       learn
                                   pandas
                          TensorFlow
                        matpletlib O PyTorch
def signal(context):
  test_code = """import torch
  list_of_tensors = [torch.randn(3),
     torch.randn(3),
     torch.randn(3)]
  try:
    exec(test_code + output)
     return 1
  except Exception as e:
     return 0
```

- Code:
 - Static analysis
 - Run tests check for errors
- Robotics: Simulate a plan and check goal



<u>Goal</u>

- Code:
 - Static analysis
 - Run tests check for errors
- Robotics: Simulate a plan and check goal
- Chemistry: Compute molecular properties





- Code:
 - Static analysis
 - Run tests check for errors
- Robotics: Simulate a plan and check goal
- Chemistry: Compute molecular properties

- Code:
 - Static analysis
 - Run tests check for errors
- Robotics: Simulate a plan and check goal
- Chemistry: Compute molecular properties

continuous/binary

- Code:
 - Static analysis
 - Run tests check for errors
- Robotics: Simulate a plan and check goal
- Chemistry: Compute molecular properties

continuous/binary

cheap/expensive

- Code:
 - Static analysis
 - Run tests check for errors
- Robotics: Simulate a plan and check goal
- Chemistry: Compute molecular properties

continuous/binary

cheap/expensive

token-by-token/sparse





Key idea: controlled generation as inference



Key idea: controlled generation as inference LM takes in a prompt

Prompt: How do I stack a list "tensors" of pytorch tensors?

Slide adapted from Alex Lew



Key idea: controlled generation as inference LM takes in a prompt, acts as a prior p_{LM}

Prompt: How do I stack a list "tensors" of pytorch tensors?

 $p_{LM}(x)$

Slide adapted from Alex Lew

torch.tensor(tensors) tensors for a list of tensors, we have to... in python, we can stack... sure, let me see... torch.stack(*tensors) torch.stack(tensors) torch.tensor(*tensors, axis=1) np.stack(tensors, axis=0)

torch.stack(tensors, dim=0)

all strings x



Key idea: controlled generation as inference LM takes in a prompt, acts as a prior p_{IM} , and ϕ s as likelihoods

Prompt: How do I stack a list "tensors" of pytorch tensors?



Slide adapted from Alex Lew

Key idea: controlled generation as inference LM takes in a prompt, acts as a *prior* p_{LM} , and ϕ s as likelihoods

Prompt: How do I stack a list "tensors" of pytorch tensors?

 $\propto p_{IM}(x)\phi_{linter}(x)\phi_{test}(x)$ def Φ_test(context): test_code = """import torch list_of_tensors = [torch.randn(3), torch.randn(3), torch.randn(3)]""" try: exec(test_code + output) return 1 except Exception as e: SisciPy learn return 0 pandas **TensorFlow** matp <a>Tib O PyTorch Slide adapted from Alex Lew

?

 $\phi_{test} = 0$ $\phi_{test} = 0$ torch.tensor(tensors) tensors $\phi_{linter} = 0$ for a list of tensors, we have to... $\phi_{linter} = 0$ in python, we can stack... $\phi_{linter} = 0$ sure, let me see... torch.stack(*tensors) torch.stack(tensors) $\phi_{test} = 0$ $\phi_{test} = 0$ torch.tensor(*tensors, axis=1) np.stack(tensors, axis=0)

torch.stack(tensors, dim=0)

all strings *x*



| Sample-rank (best-of-N, rejection sampling etc.) | |
|---|--|
| | |
| | |
| | |

e.g. best-of-n with a reward model (Nakano et al., 2021; Krishna et al., 2022; Zhou et al., 2023; Gui et al., 2024; Mudgal et al., 2024; Ichihara et al., 2025), filtering with a verifier (Olausson et al., 2023; Chen et al. 2024; Lightman et al., 2024; Xin et al., 2024)



Generations

- a = torch.stack((tensors)
- a = torch.tensor(tensors)
- a = torch.stack(tensors)

e.g. best-of-n with a reward model (Nakano et al., 2021; Krishna et al., 2022; Zhou et al., 2023; Gui et al., 2024; Mudgal et al., 2024; Ichihara et al., 2025), filtering with a verifier (Olausson et al., 2023; Chen et al. 2024; Lightman et al., 2024; Xin et al., 2024)





e.g. best-of-n with a reward model (Nakano et al., 2021; Krishna et al., 2022; Zhou et al., 2023; Gui et al., 2024; Mudgal et al., 2024; Ichihara et al., 2025), filtering with a verifier (Olausson et al., 2023; Chen et al. 2024; Lightman et al., 2024; Xin et al., 2024)

Potential Scores

SyntaxError: mismatched parentheses

ValueError: only one element tensors can be converted to Python scalars





e.g. best-of-n with a reward model (Nakano et al., 2021; Krishna et al., 2022; Zhou et al., 2023; Gui et al., 2024; Mudgal et al., 2024; Ichihara et al., 2025), filtering with a verifier (Olausson et al., 2023; Chen et al. 2024; Lightman et al., 2024; Xin et al., 2024)







Samples from the right distribution?

How many potential calls?





Distribution we sample from

 $p_{IM}(x)$



Distribution we sample from

 $p_{IM}(x)\phi(x)$

Function we score by





Distribution we sample from

Function we score by







Sample-rank needs a lot of LM calls

 $g(x) \propto p_{IM}(x)\phi(x)$

Sample-rank needs a lot of LM calls

 $g(x) \propto p_{IM}(x)\phi(x)$

Required # LM samples: $e^{D_{KL}(g||p_{LM})}$

(Chatterjee and Diaconis, 2018)





Popular approach 2: Locally Constrained Decoding

(e.g., OpenAl 2024; Shin et al., 2021; Scholak et al., 2021; Poesia et al., 2022; Willard & Louf, 2023; Moskal et al., 2024; Ugare et al., 2024)


Popular approach 2: Locally Constrained Decoding At each step:

a = torch.



Popular approach 2: Locally Constrained Decoding At each step: get next token probs





Popular approach 2: Locally Constrained Decoding At each step: get next token probs, multiply potentials





Popular approach 2: Locally Constrained Decoding At each step: get next token probs, multiply potentials, normalize





Popular approach 2: Locally Constrained Decoding At each step: get next token probs, multiply potentials, normalize, sample





Popular approach 2: Locally Constrained Decoding At each step: get next token probs, multiply potentials, normalize, sample





Popular approach 2: Locally Constrained Decoding At each step: get next token probs, multiply potentials, normalize, sample





Popular approach 2: Locally Constrained Decoding

At each step: get next token probs, multiply potentials, normalize, sample











Locally Constrained Decoding calls potentials ~100k times per step





























Locally Constrained Decoding samples from the wrong distribution Distribution has the right *support*, but overrepresents greedy samples

if tensors in torch.tensor...

for i in torch.stack(tensors):...



with torch.cuda.device('0'):...

Locally Constrained Decoding samples from the wrong distribution Distribution has the right *support*, but overrepresents greedy samples

if tensors in torch.tensor...

for i in torch.stack(tensors):...

$$l_{\phi}(x) = \prod_{i}^{|x|} \frac{p_{L}}{\sum_{x'}}$$



 $M(x_t | x_{< t}) \phi(x_t | x_{< t})$ $p_{LM}(x' | x_{< t}) \phi(x' | x_{< t})$



Better approaches to controlled generation

| | Sample-rar |
|---------------------|------------------------------|
| Output distribution | $g(x) \propto p_{LM}(x)q$ |
| # potential calls | # LM sample |
| # LM samples | $\approx \exp D_{KL}(g)$ |



Better approaches to controlled generation

| | Sample-rar |
|---------------------|-----------------------------|
| Output distribution | $g(x) \propto p_{LM}(x)q$ |
| # potential calls | # LM sample |
| # LM samples | $\approx \exp D_{KL}(g $ |



Correcting the output distribution









$$l_{\phi}(x) \prod_{i}^{|x|} w_{t}(x)$$

 $p(x) = p_{LM}(x)\phi(x)$

Better approaches to controlled generation

| | Sample-rar |
|---------------------|-----------------------------|
| Output distribution | $g(x) \propto p_{LM}(x)q$ |
| # potential calls | # LM sample |
| # LM samples | $\approx \exp D_{KL}(g $ |


| | Sample-rar |
|---------------------|----------------------------|
| Output distribution | $g(x) \propto p_{LM}(x)q$ |
| # potential calls | # LM sample |
| # LM samples | $\approx \exp D_{KL}(g)$ |



 $\Phi_{\rm eff}$: grammars, finite-state machines, regular expressions, reward models...

 $\Phi_{\rm eff}$: grammars, finite-state machines, regular expressions, reward models...

 Φ_{exp} : executing test cases, running simulations, LM as judge...

 $\Phi_{\rm eff}$: grammars, finite-state machines, regular expressions, reward models...

 Φ_{exp} : executing test cases, running simulations, LM as judge...

Sample **x** from Locally Constrained Decoding using only $\Phi_{\!\scriptscriptstyle eff}$

 $\Phi_{\rm eff}$: grammars, finite-state machines, regular expressions, reward models...

 Φ_{exp} : executing test cases, running simulations, LM as judge...

Sample **x** from Locally Constrained Decoding using only $\Phi_{\rm eff}$



| | Sample-rar |
|---------------------|----------------------------|
| Output distribution | $g(x) \propto p_{LM}(x)q$ |
| # potential calls | # LM sample |
| # LM samples | $\approx \exp D_{KL}(g)$ |



| | Sample-rar |
|---------------------|------------------------------|
| Output distribution | $g(x) \propto p_{LM}(x)q$ |
| # potential calls | # LM sample |
| # LM samples | $\approx \exp D_{KL}(g)$ |



| | Sample-rank | Locally-Constrained Decoding | Syntactic and Semantic Importance Sampling (Ours) |
|---------------------|---|---|--|
| Output distribution | $g(x) \propto p_{LM}(x)\phi(x)$ | $l_{\Phi}(x) = \prod_{t=1}^{ x } \frac{p(x_t \mathbf{x}_{< t})\phi(x_t \mathbf{x}_{< t})}{\sum_{x'} p(x' \mathbf{x}_{< t})\phi(x' \mathbf{x}_{< t})}$ | $g(x) \propto p_{LM}(x)\phi(x)$ |
| # potential calls | # LM samples | $ vocabulary \times x $ | $ \Phi_{exp}: # LM samples \Phi_{eff}: # LM samples × vocabulary × x $ |
| # LM samples | $\approx \exp D_{KL}(g \mid \mid p_{LM})$ | 1 | |

| | Sample-rank | Locally-Constrained Decoding | Syntactic and Semantic Importance Sampling (Ours) |
|---------------------|---|---|--|
| Output distribution | $g(x) \propto p_{LM}(x)\phi(x)$ | $l_{\Phi}(x) = \prod_{t=1}^{ x } \frac{p(x_t \mathbf{x}_{< t})\phi(x_t \mathbf{x}_{< t})}{\sum_{x'} p(x' \mathbf{x}_{< t})\phi(x' \mathbf{x}_{< t})}$ | $g(x) \propto p_{LM}(x)\phi(x)$ |
| # potential calls | # LM samples | $ vocabulary \times x $ | $\Phi_{exp}: # LM samples$ $\Phi_{eff}: # LM samples \times$ $ vocabulary \times x $ |
| # LM samples | $\approx \exp D_{KL}(g \mid \mid p_{LM})$ | 1 | $\approx \exp D_{KL}(g \mid \mid l_{\Phi})$ |













Correct for greediness and for Φ_{exp}

Incremental weights W_{t}

0.01

0.32

0.32







More efficient approaches to controlled generation



| ly-Constrained Decoding | Syntactic and SemanticSyntactic an SemanticImportanceSequential Mo Carlo (Ours) | |
|---|--|---|
| $\prod_{t=1}^{ x } \frac{p(x_t \mathbf{x}_{< t})\phi(x_t \mathbf{x}_{< t})}{\sum_{x'} p(x' \mathbf{x}_{< t})\phi(x' \mathbf{x}_{< t})}$ | $g(x) \propto p_{LM}(x)\phi(x)$ | $g(x) \propto p_{LM}(x)\phi(x)$ |
| | <pre></pre> | Φ_{exp} # LM samples × x |
| abulary $ \times \mathbf{x} $ | Φ_{eff} : # LM samples × vocabulary × x 📀 | Φ_{eff} : # LM samples × vocabulary × x 📀 |
| 1 | $\approx \exp D_{KL}(g \mid \mid l_{\Phi})$ | ~10x fewer than Syntactic and Semantic Importance Sampling |

More efficient approaches to controlled generation How do these methods perform?

| | Sample-rank | Locally-Constrained Decoding | Syntactic and Semantic Importance Sampling (Ours) | Syntactic and Semantic Sequential Monte Carlo (Ours) |
|---------------------|------------------------------------|---|--|---|
| Output distribution | $g(x) \propto p_{LM}(x)\phi(x)$ | $l_{\Phi}(x) = \prod_{t=1}^{ x } \frac{p(x_t \mathbf{x}_{< t})\phi(x_t \mathbf{x}_{< t})}{\sum_{x'} p(x' \mathbf{x}_{< t})\phi(x' \mathbf{x}_{< t})}$ | $g(x) \propto p_{LM}(x)\phi(x)$ | $g(x) \propto p_{LM}(x)\phi(x)$ |
| # verifier calls | # LM samples | vocabulary × x | $\Phi_{exp}: # LM samples$ $\Phi_{eff}: # LM samples \times$ $ vocabulary \times x $ | $\Phi_{exp} : # LM samples \times x $ $\Phi_{eff} : # LM samples \times x $ $ vocabulary \times x $ |
| # LM samples | $\approx \exp D_{KL}(g p_{LM})$ | 1 | $\approx \exp D_{KL}(g \mid \mid l_{\Phi})$ | ~10x fewer than Syntactic and Semantic Importance Sampling |

Sequential Monte Carlo boosts performance across challenging domains



Sequential Monte Carlo is more frugal than Importance Sampling ~10x fewer LM calls in domains with constraining potentials



LM samples



Sequential Monte Carlo allows small LMs to punch above their weights





+ greediness correction, expensive potentials + SMC



tl;dr

Frame controlled generation as probabilistic inference Use Sequential Monte Carlo to sample from the posterior

Library: <u>github.com/genlm/genlm-control</u>

Our org GenLM is recruiting! genlm.org



| <pre># Create a language model potential. llm = PromptedLLM.from_name("gpt2") llm.set_prompt_from_str("Here is my honest opinion:")</pre> |
|---|
| # Create a finite-state automaton potential using a regular expression. fsa = BoolFSA.from_regex(r" SMC is (🍦 🍦 🥶 🥶 🌛 🌛) with LMs") |
| # Coerce the FSA so that it operates on the token type of the language model. coerced_fsa = fsa.coerce(llm, f=b"".join) |
| <pre># Create a token sampler that combines the language model and FSA. token_sampler = AWRS(llm, coerced_fsa)</pre> |
| <pre># Generate text using SMC. # Generation is asynchronous; use `await` if calling in an async context (like in an a # function or in a Jupyter notebook) and `asyncio.run(token_sampler.smc())` otherwi sequences = await token_sampler.smc(n_particles=10, # Number of candidate sequences to maintain ess_threshold=0.5, # Threshold for resampling max_tokens=30, # Maximum sequence length verbosity=1 # Print particles at each step)</pre> |
| <pre>sequences.decoded_posterior # Example output: # { # ' SMC is 🔥 🍐 with LMs': 1.0, # }</pre> |
| |

enlm.control import PromptedLLM, BoolFSA, AWRS





Appendix

Does it matter that our output follows the right distribution? Yes! Methods that are closer in KL to true posterior perform better



+ greediness correction + Sequential Monte Carlo

- Locally Constrained Decoding + expensive potentials



Does it matter that our output follows the right distribution?



+ greediness correction

Yes! Better methods' probabilities more correlated with downstream performance

+ expensive potentials

+ Sequential Monte Carlo

