

# Speed-Accuracy Tradeoffs in Tagging with Variable-Order CRFs and Structured Sparsity

Tim Vieira,\* Ryan Cotterell\* and Jason Eisner

Johns Hopkins University

## Variable-Order CRFs

**Goal:** Define a good conditional distribution over tag sequences.

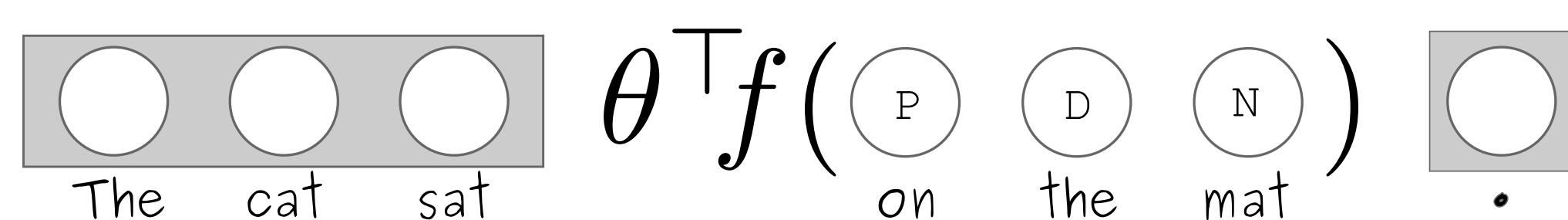
$$p_{\theta} \left( \begin{array}{cccccc} \text{D} & \text{N} & \text{V} & \text{P} & \text{D} & \text{N} & \cdot \\ \text{The} & \text{cat} & \text{sat} & \text{on} & \text{the} & \text{mat} & \cdot \end{array} \mid \mathcal{X} \right)$$

Certain combinations go well together, some don't.

label-word: D-the ✓, v-the ✗

label-label: D-N ✓, D-V ✗

Sometimes, it's useful to look at larger label combinations.



**The problem:** For features to look at output contexts of size  $k$ , we need  $\mathcal{O}(n \cdot |Y|^k)$  time for inference even if most combinations don't improve the model, e.g., combinations that are easily ruled out by local features.

$$p_{\theta}(y|x) = \frac{1}{Z_{\theta}(x)} \exp \left( \sum_{t=1}^{n+1} \theta^{\top} f(x, t, y_{t-k-1} \dots y_t) \right)$$

**The VoCRF idea:** Remove output contexts that aren't necessary!

Turn  $\mathcal{O}(n \cdot |Y|^k)$  into  $\mathcal{O}(n \cdot |\bar{\mathcal{W}}|)$

all contexts → important contexts

**Technical details:**

Need closure under prefixes & last-character substitution

$$\mathcal{W} \rightarrow \bar{\mathcal{W}}$$

Even if you don't use those features

Can lift assumption with  $\phi$ -arcs.

**Very flexible**

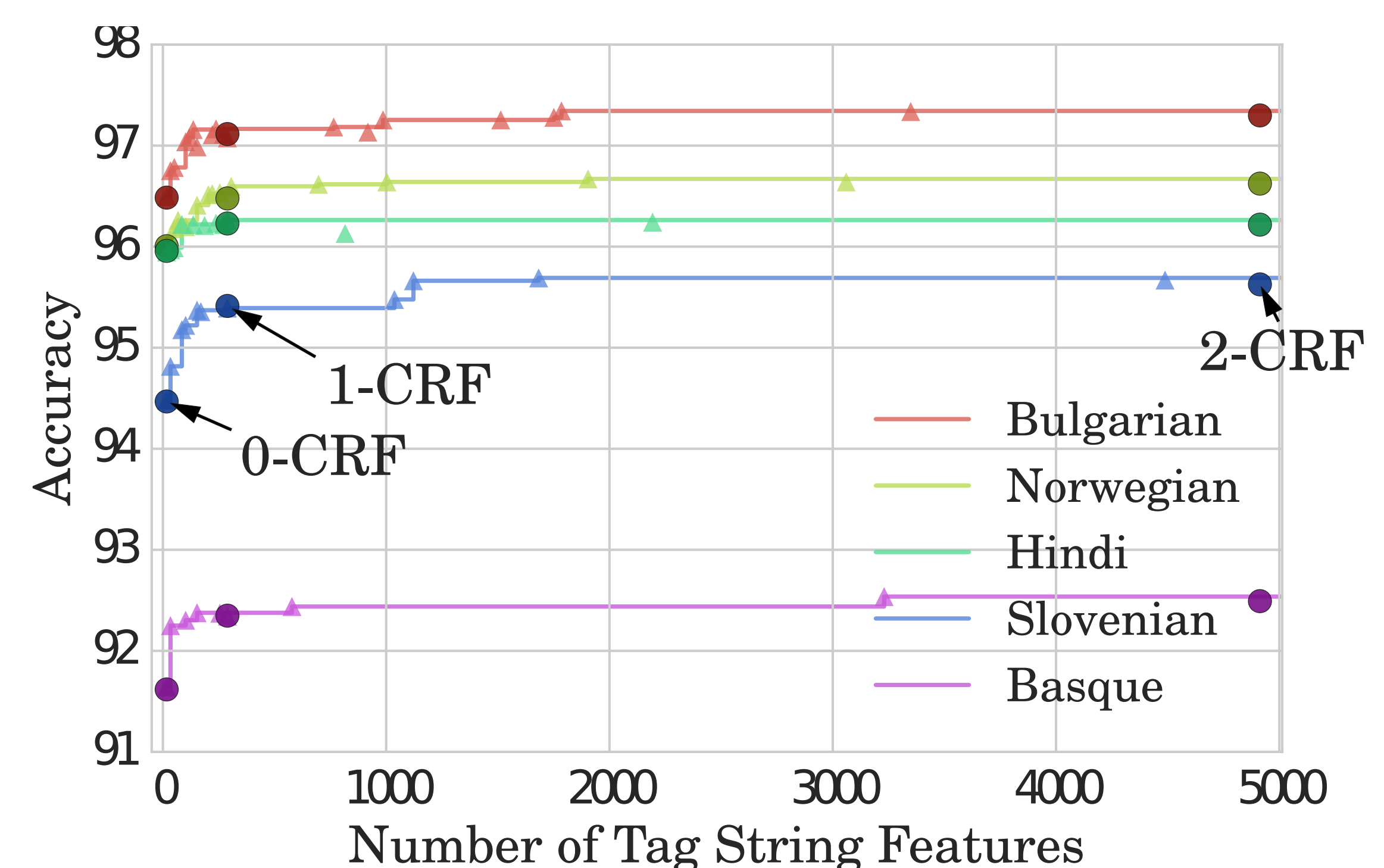
- No need to specify a fixed size.
- Covers semi-Markov & higher-order
- (Ye et al., 2009, Cuong et al., 2014)

- Can use any subset of  $Y^*$ .
- Easy to implement!
- One alg. many models.

**"Correcting" prior work**

Original algorithm for computing gradients and expectations was unnecessarily slow and complicated. Our revised algorithm is  $\mathcal{O}(|\bar{\mathcal{W}}|)$  times faster  $\mathcal{O}(\text{a few pages})$  simpler!

- Just run autodiff on their forward algorithm!
- Protip: Evaluating the gradient should be as fast as the function!
- Check out Jason's paper at the structured prediction workshop for more on the connection between autodiff and inference.



## Structured Sparsity

**Goal:** Select higher-order features  $\mathcal{W}$ , which gives us the best possible accuracy under a budget for runtime.

**How:** Augmenting the training objective with a penalty for runtime!

$$\sum_{i=1}^m -\log p_{\theta}(y^{(i)} | x^{(i)}) + \lambda \|\theta\|_2^2 + \gamma \mathcal{R}(\theta)$$

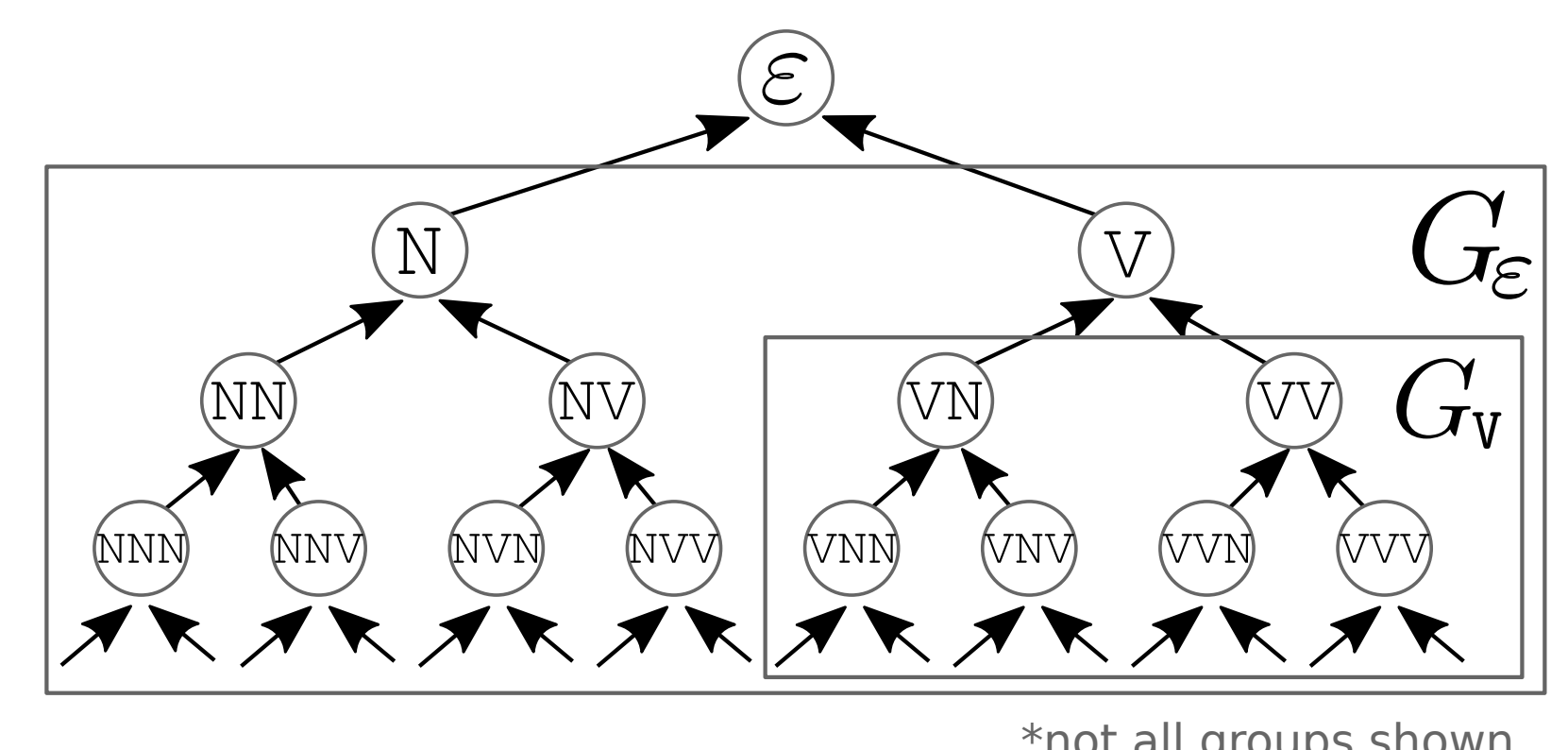
loss                      regularizer                      runtime

$\theta$  implicitly encodes  $\mathcal{W}$  in its nonzero entries.

Sparsity → Speed

Dependencies among features

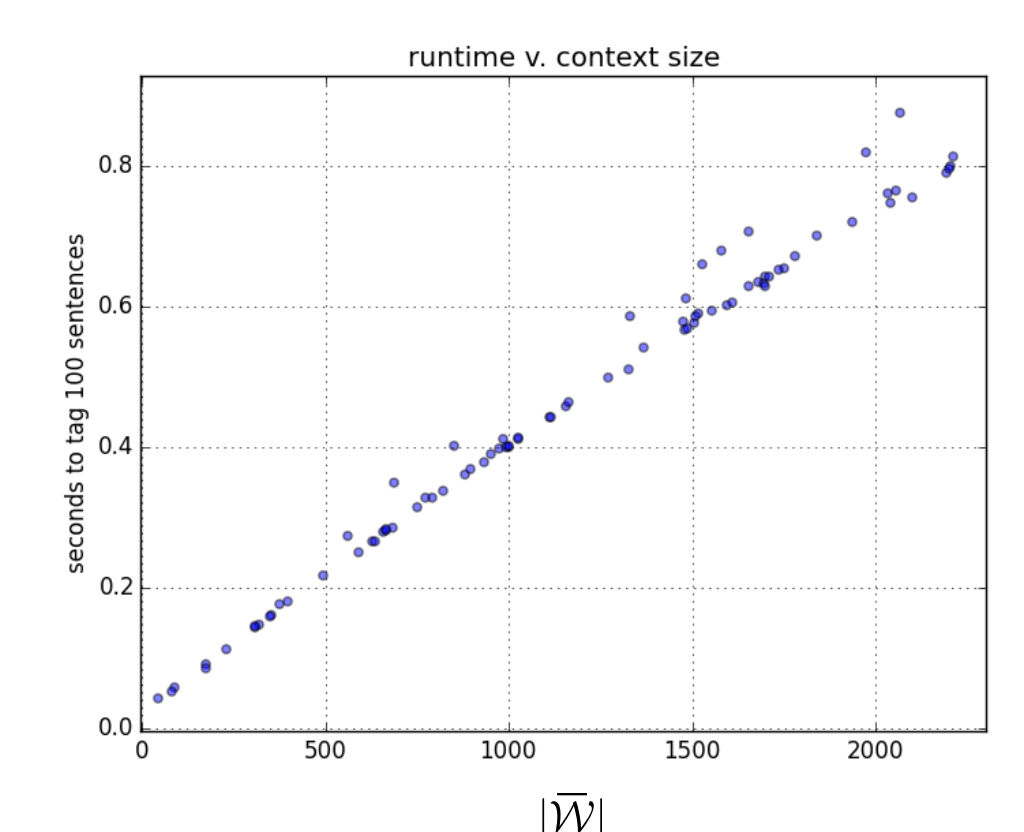
- prefix closure  
NNV → NN → N →  $\epsilon$
- last tag subst. closure  
NNV → NNV  
NN → NV



**Ideal runtime**

$$\mathcal{R}^*(\theta) = |\bar{\mathcal{W}}| = \sum_{w \in Y^*} \|\theta_{G_w}\|_0$$

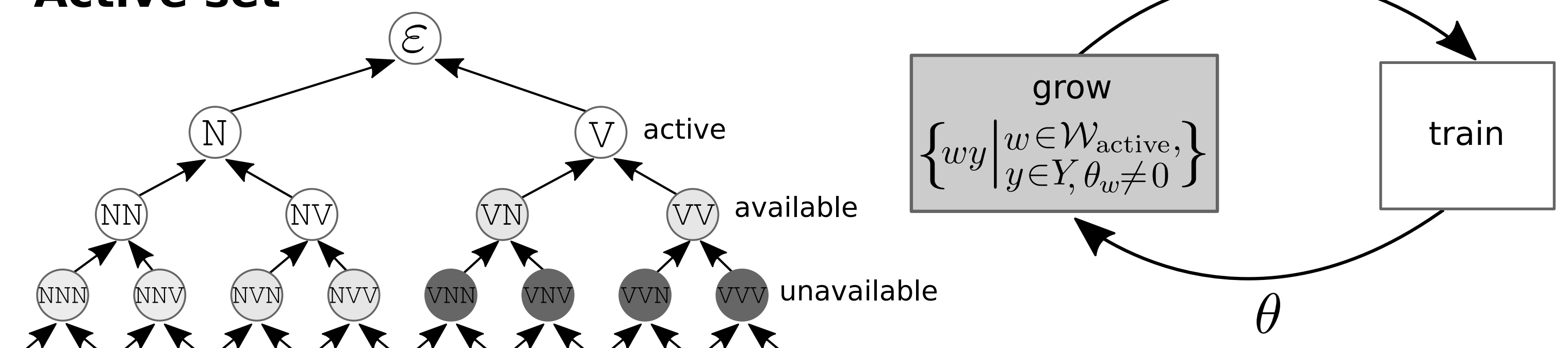
too hard to optimize!



**Convex surrogate**

$$\mathcal{R}(\theta) = \sum_{w \in Y^*} \|\theta_{G_w}\|_2 \quad \text{group lasso}$$

**Active set**



## Experiments

- Part of speech tagging with Universal Tags in 5 languages.
- Best system in **bold**.

- Superscript k indicates a significant difference from the k-CRF's accuracy (paired-permutation  $p < 0.5$ ), color indicates **better** or **worse**.

- Underlined system is the fastest "statistically indistinguishable" model compared to the 2-CRF.

	k-CRF ( $ \bar{\mathcal{W}}  = 17^{k+1}$ )			VoCRF at different model sizes $ \bar{\mathcal{W}} $ (which is proportional to runtime)									
	0 (17)	1 (289)	2 (4913)	≤ 34	≤ 85	≤ 170	≤ 340	≤ 850	≤ 1700	≤ 2550	≤ 3400	≤ 4250	≤ 5100
Ba	91.61 <sup>1,2</sup>	92.35 <sup>0</sup>	92.49 <sup>0</sup>	92.25 <sup>0,2</sup>	92.25 <sup>0,2</sup>	92.38 <sup>0</sup>	92.34 <sup>0</sup>	92.44 <sup>0</sup>	92.44 <sup>0</sup>	92.44 <sup>0</sup>	<b>92.54<sup>0</sup></b>	92.54 <sup>0</sup>	92.54 <sup>0</sup>
Bu	96.48 <sup>1,2</sup>	97.11 <sup>0,2</sup>	97.29 <sup>0,1</sup>	96.75 <sup>0,1,2</sup>	96.78 <sup>0,1,2</sup>	96.99 <sup>0,1,2</sup>	97.08 <sup>0,2</sup>	97.18 <sup>0,1</sup>	97.25 <sup>0,1</sup>	<b>97.34<sup>0,1</sup></b>	97.34 <sup>0,1</sup>	97.34 <sup>0,1</sup>	97.34 <sup>0,1</sup>
Hi	95.96 <sup>1,2</sup>	96.22 <sup>0</sup>	96.21 <sup>0</sup>	95.97 <sup>1,2</sup>	96.22 <sup>0</sup>	96.22 <sup>0</sup>	96.26 <sup>0</sup>	96.13 <sup>0</sup>	96.13 <sup>0</sup>	<b>96.24<sup>0</sup></b>	96.24 <sup>0</sup>	96.24 <sup>0</sup>	96.24 <sup>0</sup>
No	96.00 <sup>1,2</sup>	<u>96.64<sup>0</sup></u>	96.66 <sup>0</sup>	96.07 <sup>1,2</sup>	96.26 <sup>0,1,2</sup>	<u>96.41<sup>0</sup></u>	96.60 <sup>0</sup>	96.62 <sup>0</sup>	96.64 <sup>0</sup>	<b>96.67<sup>0</sup></b>	96.64 <sup>0</sup>	96.64 <sup>0</sup>	96.64 <sup>0</sup>
Sl	94.46 <sup>1,2</sup>	95.41 <sup>0,2</sup>	<u>95.62<sup>0,1</sup></u>	94.82 <sup>1,2</sup>	95.18 <sup>0,2</sup>	95.36 <sup>0,2</sup>	95.39 <sup>0,2</sup>	95.39 <sup>0,2</sup>	<b>95.69<sup>0,1</sup></b>	95.69 <sup>0,1</sup>	95.69 <sup>0,1</sup>	95.69 <sup>0,1</sup>	95.67 <sup>0,1</sup>