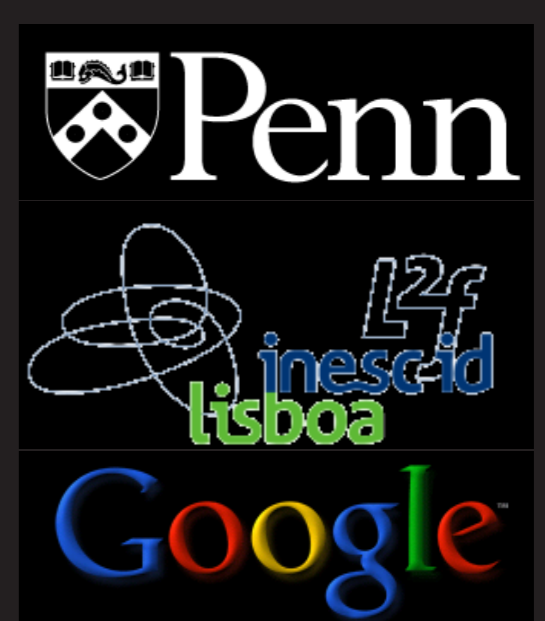


# Sparsity in Dependency Grammar Induction

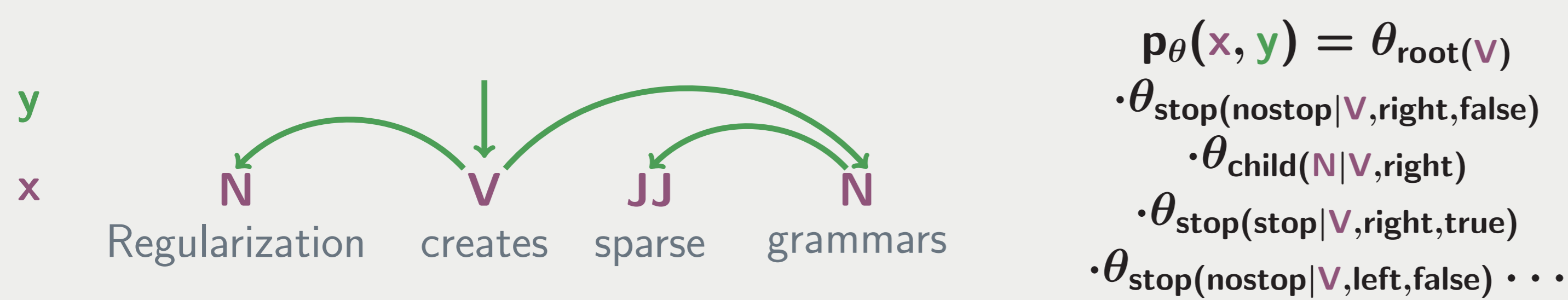
Jennifer Gillenwater<sup>1</sup>, Kuzman Ganchev<sup>1</sup>, João Graça<sup>2</sup>, Fernando Pereira<sup>3</sup>, Ben Taskar<sup>1</sup>

<sup>1</sup>University of Pennsylvania  
<sup>2</sup>L<sup>2</sup>F INESC-ID  
<sup>3</sup>Google



## Motivation

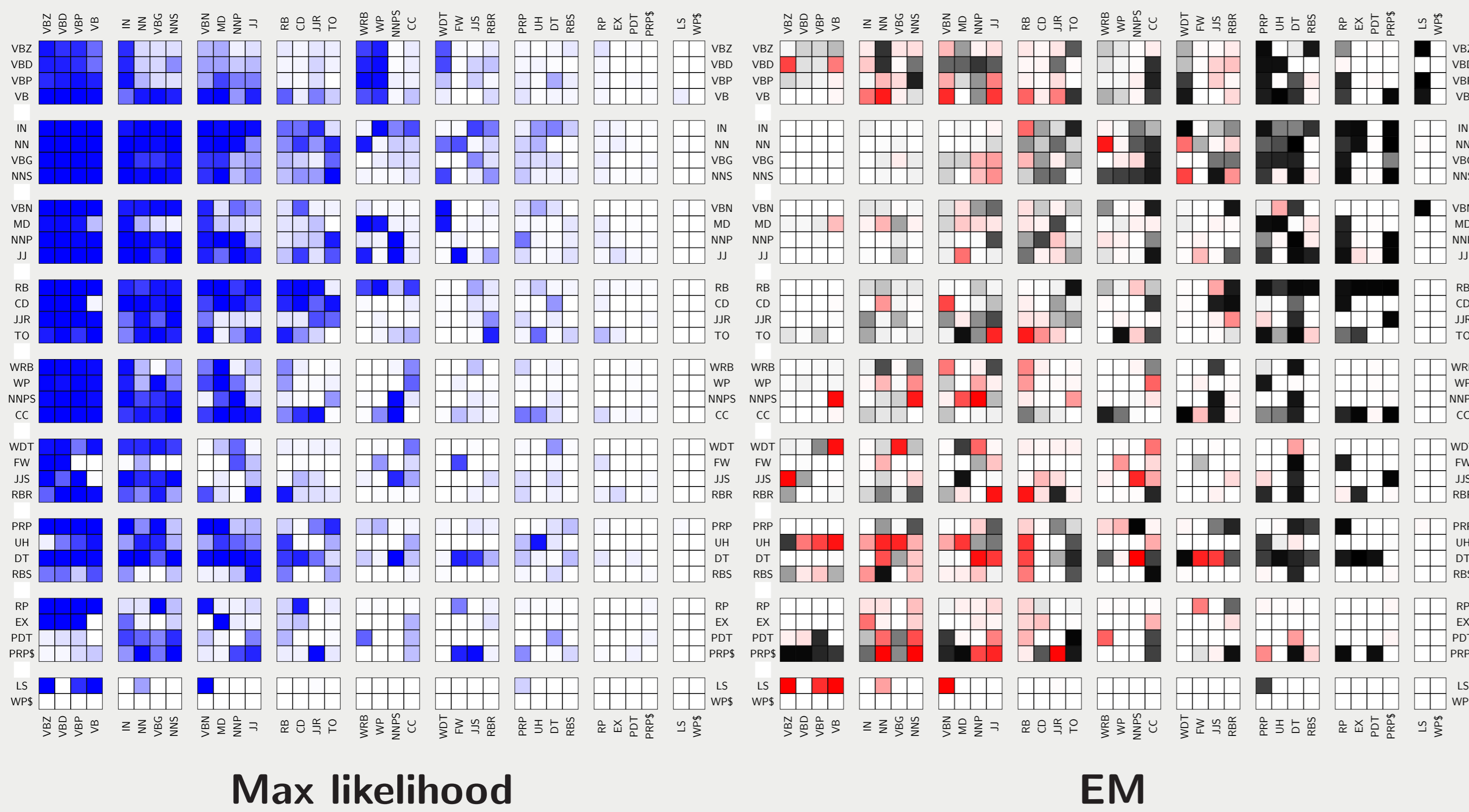
Dependency Model with Valence (Klein and Manning, ACL 2004)



- ▶ **Task:** Unsupervised dependency grammar induction
- ▶ **Problem:** Model is simple, but still **too permissive**; most relations (e.g.  $\text{DET} \rightarrow \text{V}, \text{N}, \text{JJ}, \text{etc.}$ ) should not occur
- ▶ **Solution:** **Posterior constraints** to limit grammar ambiguity during learning

## Traditional Objective Optimization

- ▶ **Traditional objective:** marginal log likelihood  $\max_{\theta} \mathcal{L}(\theta) = \hat{\mathbb{E}}_x[\log \sum_y p_{\theta}(x, y)]$
- ▶ **Optimization method:** Expectation maximization (EM)
- ▶ **Figures:** Parent tags across, child tags down
- ▶ **Left:** Blank squares have max posterior **0**; many parent-child relations don't occur
- ▶ **Right:** Red have max < supervised, black have max > supervised; many dark squares implies model assigns non-zero probability to too many pairs

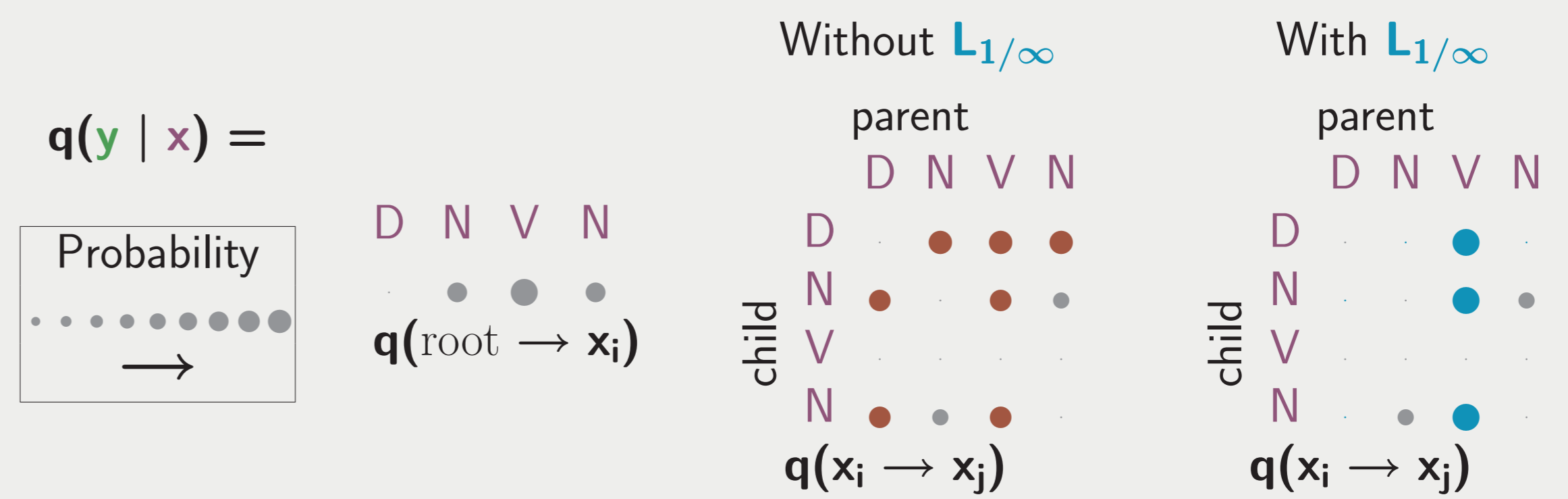


## Posterior Regularization

Minimize # of unique pairs through E-step penalty,  $L_{1/\infty}$  on the posteriors  $q(y | x)$  (Graça et al., NIPS 2007 & 2009)

$$\text{M-Step } \theta^{t+1} = \arg \max_{\theta} \hat{\mathbb{E}}_x \left[ \sum_y q^t(y | x) \log p_{\theta}(x, y) \right]$$

$$\text{E-Step } q^t(y | x) = \arg \min_{q(y|x)} \text{KL}(q(y | x) \| p_{\theta^t}(y | x)) + \sigma L_{1/\infty}(q(y | x))$$



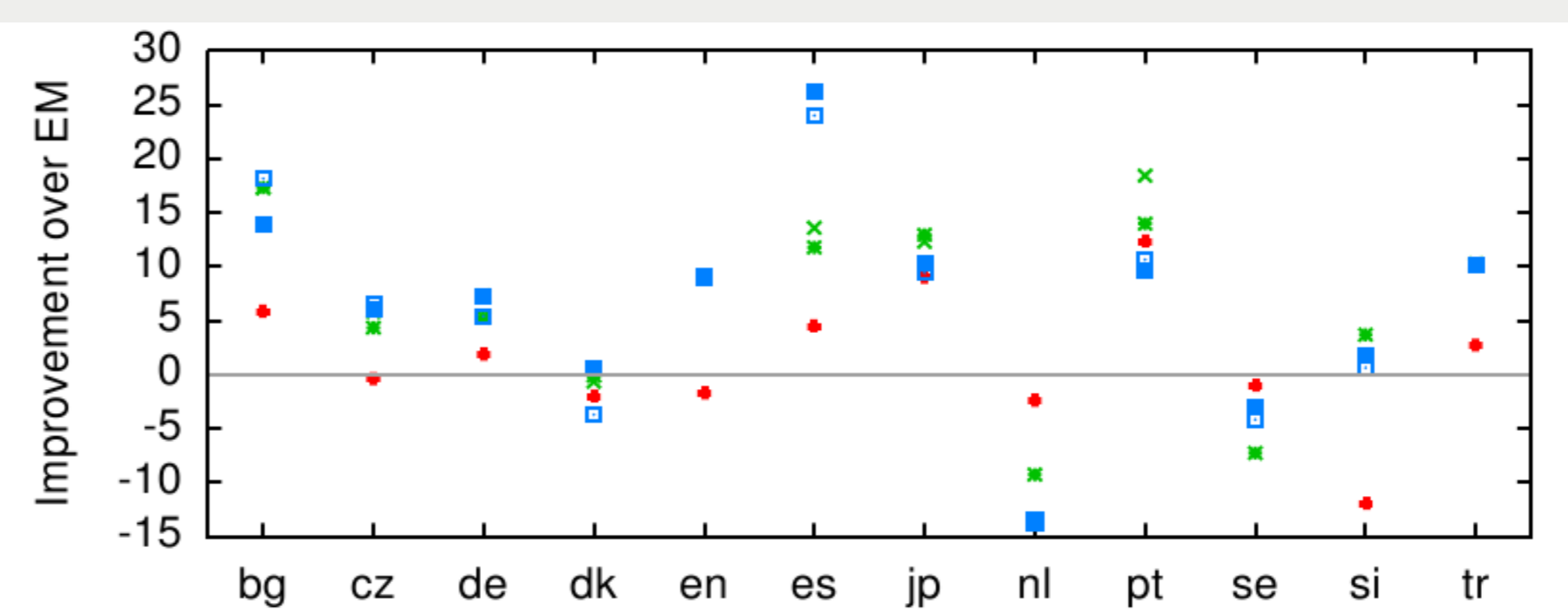
## Experiments on English

- ▶ Penn Treebank data, strip punctuation, consider separately sentences of length  $\leq 10$  and 20, initialize model "harmonically", try  $\sigma \in \{80, 100, 120, 140, 160, 180\}$
- ▶ **Table:** top — experiments on the basic dependency model with valence: our method and two parameter priors (Cohen et al., NIPS 2008; Cohen and Smith, NAACL 2009); bottom — extended version of the model with parameters that use more valence information: our method, and a non-sparsifying ( $\alpha = 1$ ) discounting Dirichlet prior (DD) with random pools initialization and learned backoff weight  $\lambda$  (Headden et al., NAACL 2009).
- ▶ PR with random pools would likely produce best result of all

Learning Method	Accuracy		
	$\leq 10$	$\leq 20$	all
PR ( $\sigma = 140$ )	<b>62.1</b>	<b>53.8</b>	<b>49.1</b>
LN families	59.3	45.1	39.0
SLN TieV & N	61.3	47.4	41.4
PR ( $\sigma = 140, \lambda = 1/3$ )	64.4	55.2	50.5
DD ( $\alpha = 1, \lambda$ learned)	<b>65.0 (<math>\pm 5.7</math>)</b>		

## Experiments on 11 Other Languages

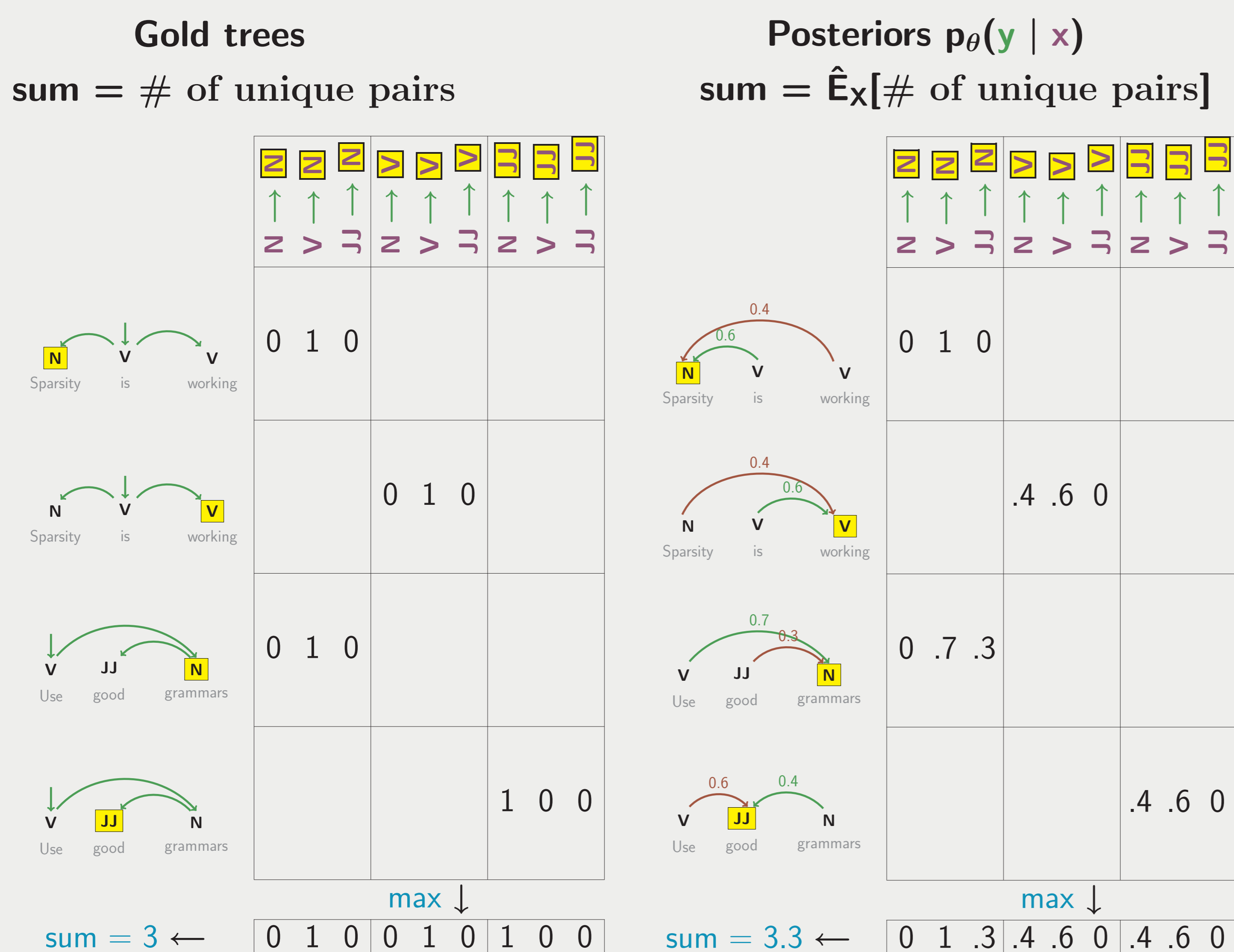
**Figure:** Relative error with respect to EM on the extended model; DD = discounting Dirichlet prior, PR = posterior regularization ( $\sigma = 160$  chosen on English), PR-S = symmetric version of constraints, PR-AS = asymmetric version; Avg = average improvement over EM, W = # of languages better than EM



## Parameter Regularization: $\mathcal{L}(\theta) + \log p(\theta)$

- ▶ Hierarchical Dirichlet processes (Liang et al., EMNLP 2007; Johnson et al., NIPS 2007)
- ▶ Discounting Dirichlet prior (Headden et al., ACL 2009)
- ▶ Logistic normal prior (Cohen et al., NIPS 2008; Cohen and Smith, NAACL 2009)
- ▶ All of these tend to **reduce unique # of children per parent**, rather than directly **reducing # of unique parent-child pairs**:  $\theta_{\text{child}(y|x)} \neq \text{posterior}(x \rightarrow y)$

## Ambiguity Measure Using Posteriors: $L_{1/\infty}$



## Parse Analysis

Parse	Unique parent-child pairs
	(v, nc); (nc, d)
	(v, d); (d, nc)
	(v, nc); (v, v)

- ▶ Parses 1 and 3: 3 unique pairs total
- ▶ Parses 2 and 3: 4 unique pairs total

## Conclusion

- ▶ For the basic model, average improvements over EM are 1.6% for DD, 6.7% for PR
- ▶ For the extended model, average improvements over EM are 1.4% for DD, 6.4% for PR
- ▶ Using posterior regularization significantly improves parsing accuracy