

Graph-Based Posterior Regularization for Semi-Supervised Structured Prediction



Luheng He

Jennifer Gillenwater

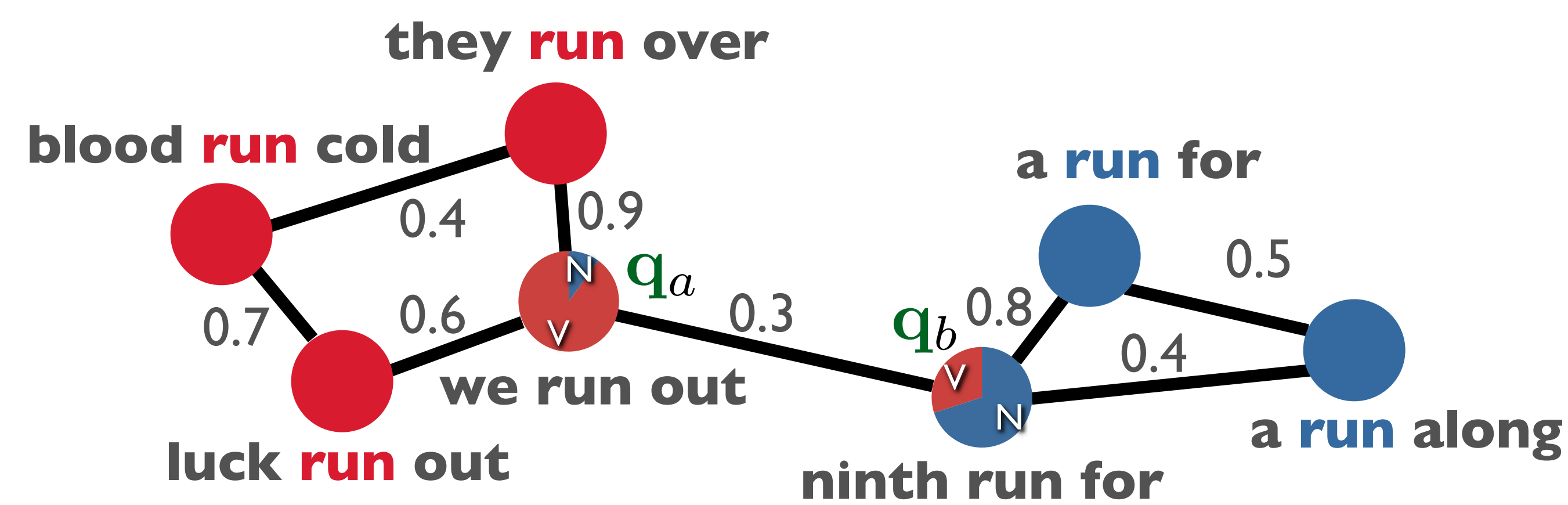
Ben Taskar

{luhe,jengi}@cis.upenn.edu, taskar@cs.washington.edu



GRAPH-BASED LEARNING

Labels: **verb (V)**, **noun (N)**, etc.

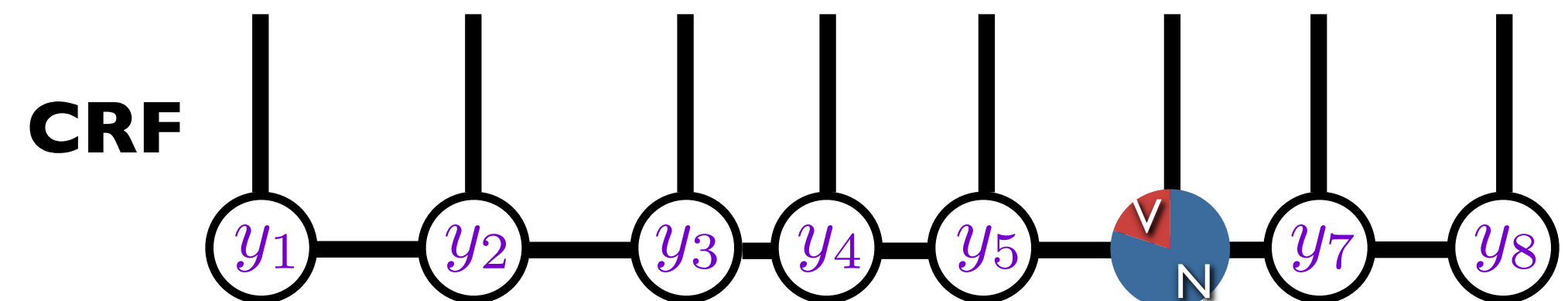


Minimize Laplacian-based objective, summing over all neighbors of unlabeled nodes:

$$\text{Lap}(q) = \sum_{a=1}^N \sum_{b=L+1}^N w_{ab} \|q_a - q_b\|^2$$

STRUCTURED PREDICTION

$x =$ **The soldiers of the ninth run for cover**



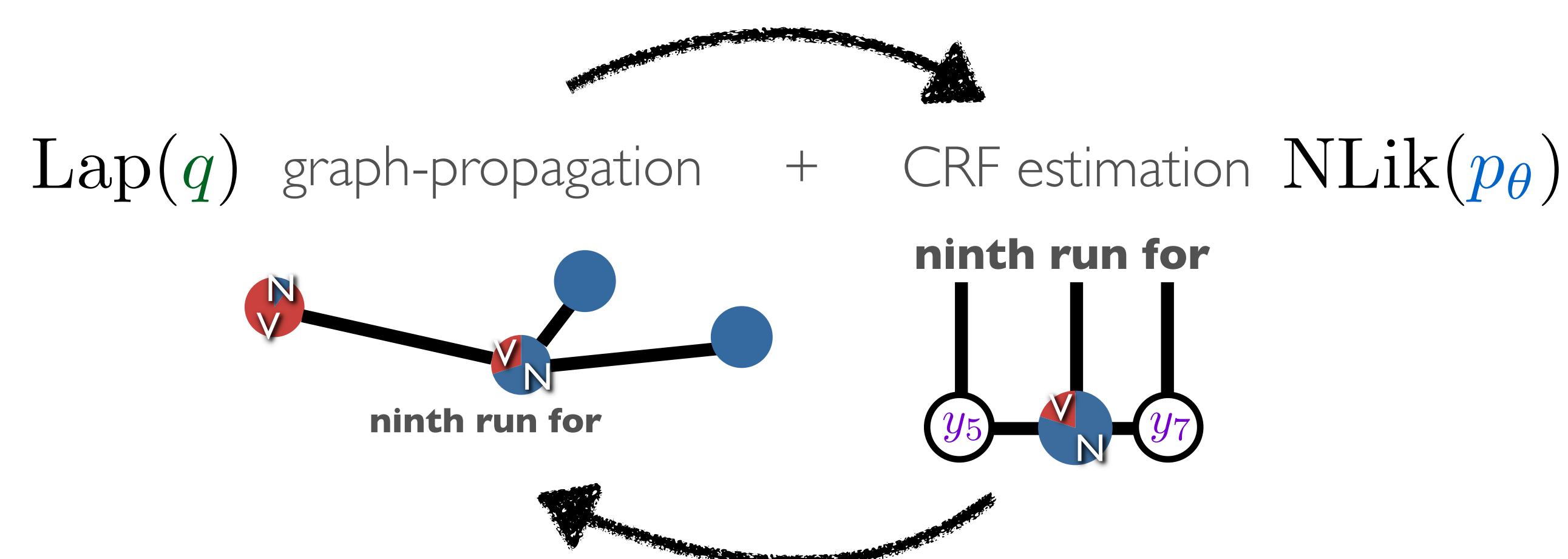
$$p_\theta(y | x) = \frac{1}{Z_\theta(x)} \exp \left[\sum_{t=1}^T \theta^\top f(y_t, y_{t-1}, x) \right]$$

Minimize negative log-likelihood, summing over all labeled sentences:

$$\text{NLik}(p_\theta) = - \sum_{i=1}^{\ell} \log p_\theta(y^i | x^i)$$

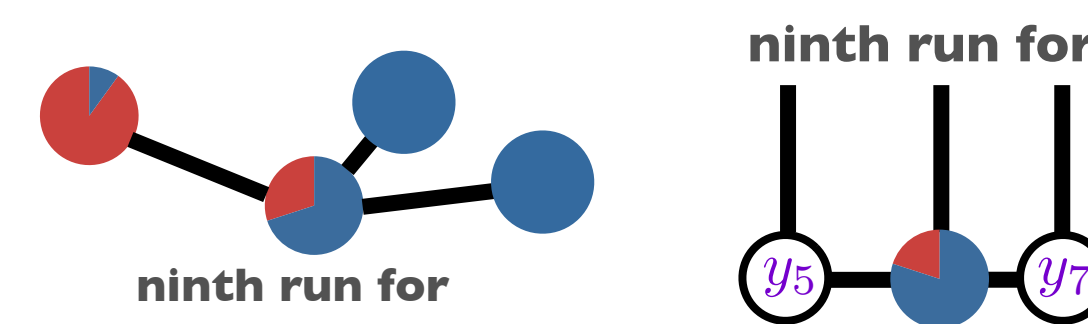
COMBINATION

Most closely related work: Subramanya et al. (EMNLP 2010) --- Iterative procedure, marginals of CRF initialize graph-propagation (GP), then GP results provide additional training data for CRF learning.



This work: retains efficiency of Subramanya et al (EMNLP 2010) while optimizing an extendible, joint objective.

JOINT OBJECTIVE



$$\mathcal{J}(q, p_\theta) = \text{Lap}(q) + \text{NLik}(p_\theta) + \text{KL}(q \| p_\theta)$$

Couple the methods via KL divergence.

(# tags)⁸ values, compactly represented by θ in the case of p

q	p_θ	The soldiers of the ninth run for cover							
7e-5	2e-5	N	N	N	N	N	N	N	N
3e-6	8e-6	N	N	N	N	N	N	N	V
...

OPTIMIZATION

p 's parameterization makes its update simple:

$$\theta \text{ update: } \theta' = \theta - \eta \frac{\partial \mathcal{J}(q, p_\theta)}{\partial \theta}$$

q has more freedom: $q^i \in \Delta$ of dimension (# tags)^(i's length)

$$\text{standard gradient update: } q_y^{i'} = \text{proj}_\Delta \left(q_y^i - \eta \frac{\partial \mathcal{J}(q, p_\theta)}{\partial q_y^i} \right)$$

-Problem 1: projection is hard $q^i \notin \Delta$

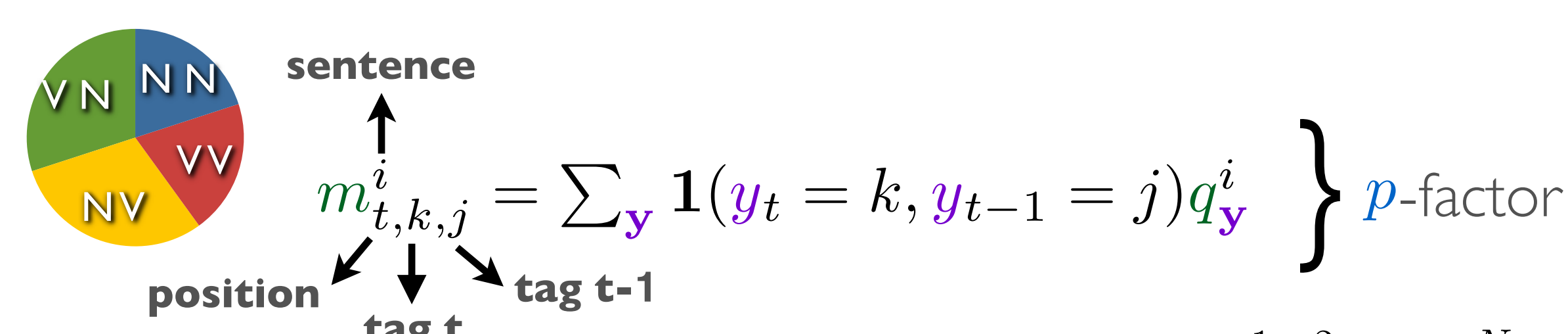
-Problem 2: no compact form (# tags)^(i's length) values

~~X~~ Standard gradient descent on the primal isn't feasible for q

What about optimizing q in the dual? $\mathcal{J}(q, p_\theta) + \gamma \left(\sum_y q_y^i - 1 \right)$

Posterior Regularization (PR) of Ganchev et al. (JMLR 2010) uses the dual, and differs from our objective only in the first term.

This work: $\text{Lap}(q) \rightarrow$ Standard PR: $\text{Linear}(m)$



$\text{Lap}(m)$ is a quadratic function though, so its dual requires an expensive matrix inverse.

$$\begin{pmatrix} 1 & & & \\ 2 & & & \\ \vdots & & & \\ N & & & \end{pmatrix}^{-1}$$

EXPONENTIATED GRADIENT

Alternative type of gradient update makes "projection" efficient:

$$q_y^{i'} = \frac{1}{Z_q(x^i)} q_y^i \exp \left[-\eta \frac{\partial \mathcal{J}(q, p_\theta)}{\partial q_y^i} \right]$$

$$\exp \left[-\eta \frac{\partial \mathcal{J}(q, p_\theta)}{\partial q_y^i} \right] =$$

$$\exp \left[-\eta \sum_{t=1}^T \frac{\partial \text{Lap}(m_y^i)}{\partial m_{t,y_t,y_{t-1}}} + \eta (\log p_\theta(y | x^i) - \log q_y^i - 1) \right]$$

$$= \exp \left[-\eta \sum_{t=1}^T \frac{\partial \text{Lap}(m_y^i)}{\partial m_{t,y_t,y_{t-1}}} \right] p_\theta(y | x^i) \eta (q_y^i)^{-\eta} e$$

product of p-factors

$\text{proj}_\Delta \rightarrow Z_q(x^i)$, computable via forward-backward

EXTENSION

$\text{Lap}(q) \rightarrow$ any convex, differentiable $g(m)$

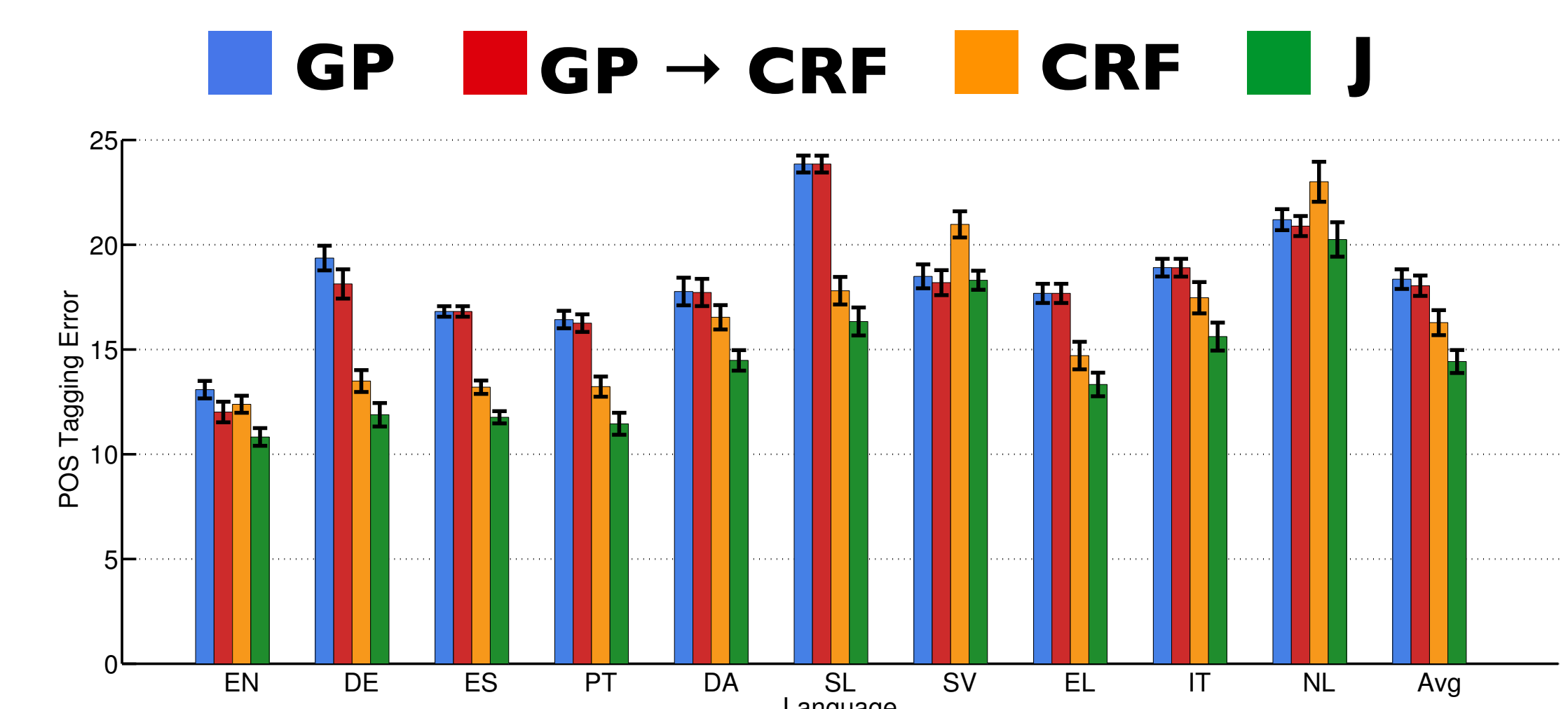
Theorem: The EM-like optimization procedure below converges to a local optimum of the joint objective

$$M\text{-step: } \theta' = \theta - \eta \frac{\partial \mathcal{J}(q, p_\theta)}{\partial \theta}$$

$$E\text{-step: } q_y^{i'} = \frac{1}{Z_q(x^i)} q_y^i \exp \left[-\eta \frac{\partial \mathcal{J}(q, p_\theta)}{\partial q_y^i} \right]$$

EXPERIMENTS

Part-of-speech tagging



Handwriting recognition

	GP	GP \rightarrow CRF	CRF	J
Mean	17.57	15.07	9.82	4.89
StdDev	0.30	0.35	0.48	0.42

Code: <https://code.google.com/p/pr-graph/>