

GRAPH-BASED
POSTERIOR REGULARIZATION
FOR
SEMI-SUPERVISED
STRUCTURED PREDICTION

Luheng He Jennifer Gillenwater
University of Pennsylvania

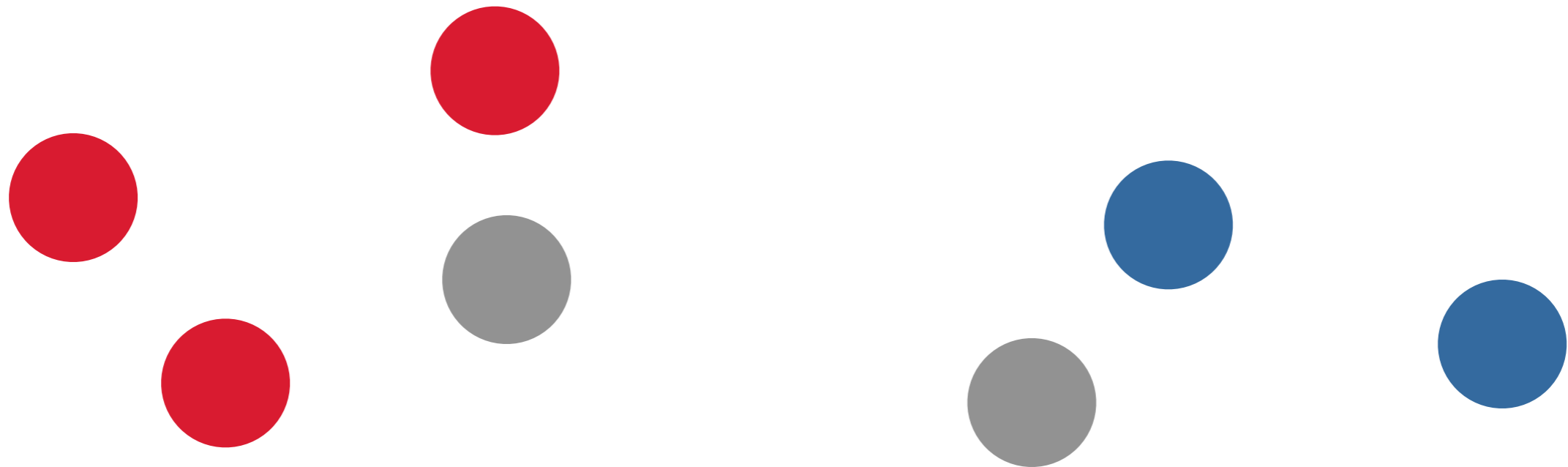
Ben Taskar
University of Washington

GRAPH-BASED LEARNING

GRAPH-BASED LEARNING

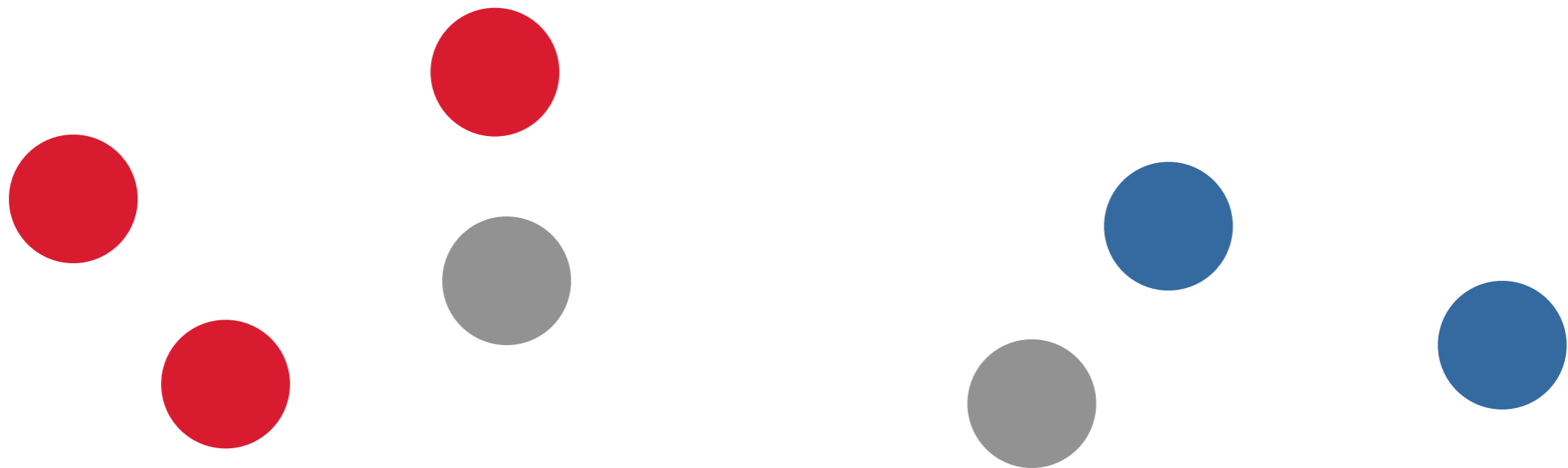


GRAPH-BASED LEARNING



GRAPH-BASED LEARNING

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



GRAPH-BASED LEARNING

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

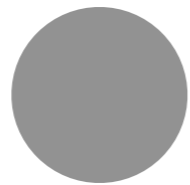
they **run** over



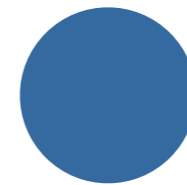
blood **run** cold



we **run** out



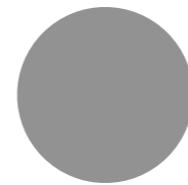
a **run** for



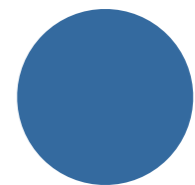
luck **run** out



ninth **run** for

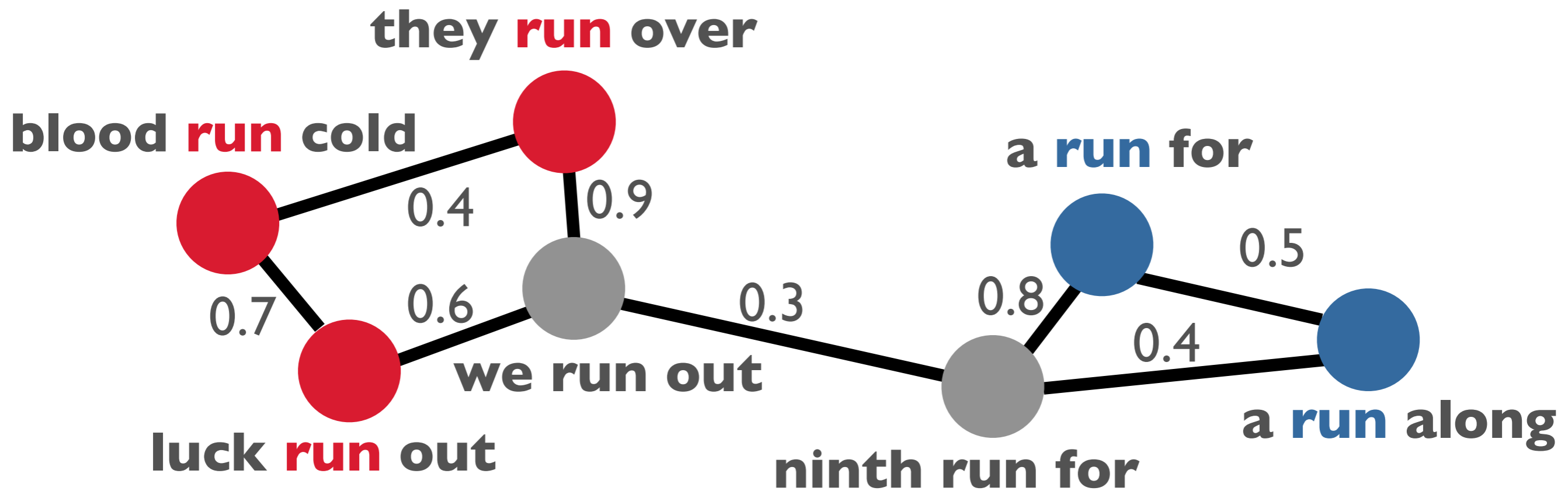


a **run** along



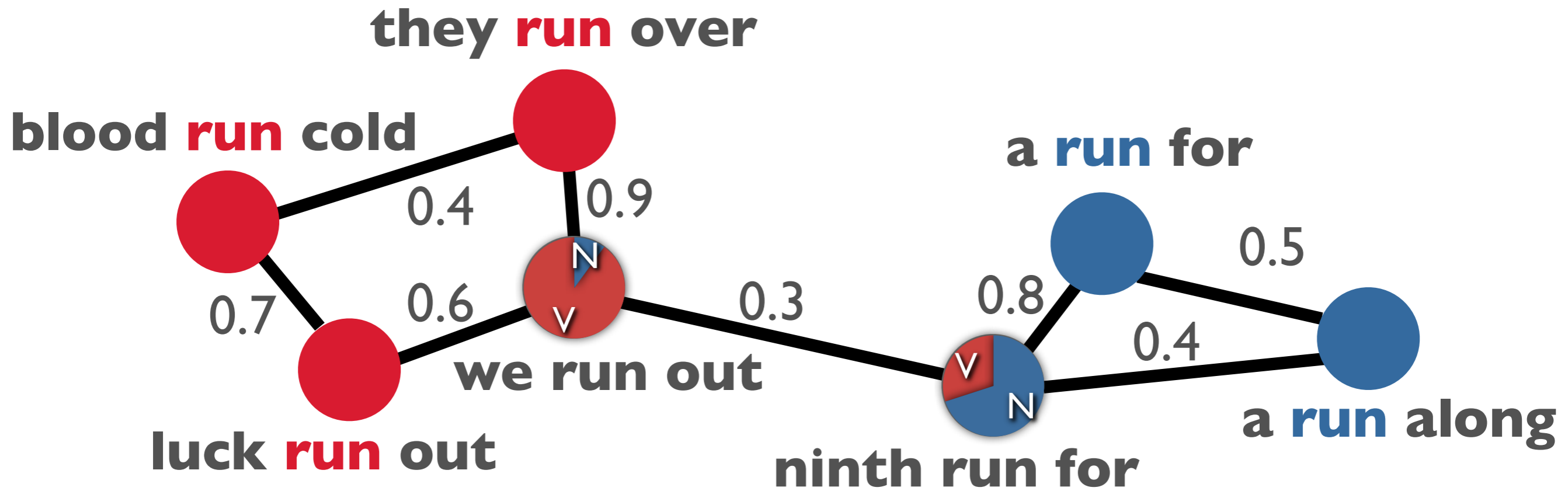
GRAPH-BASED LEARNING

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



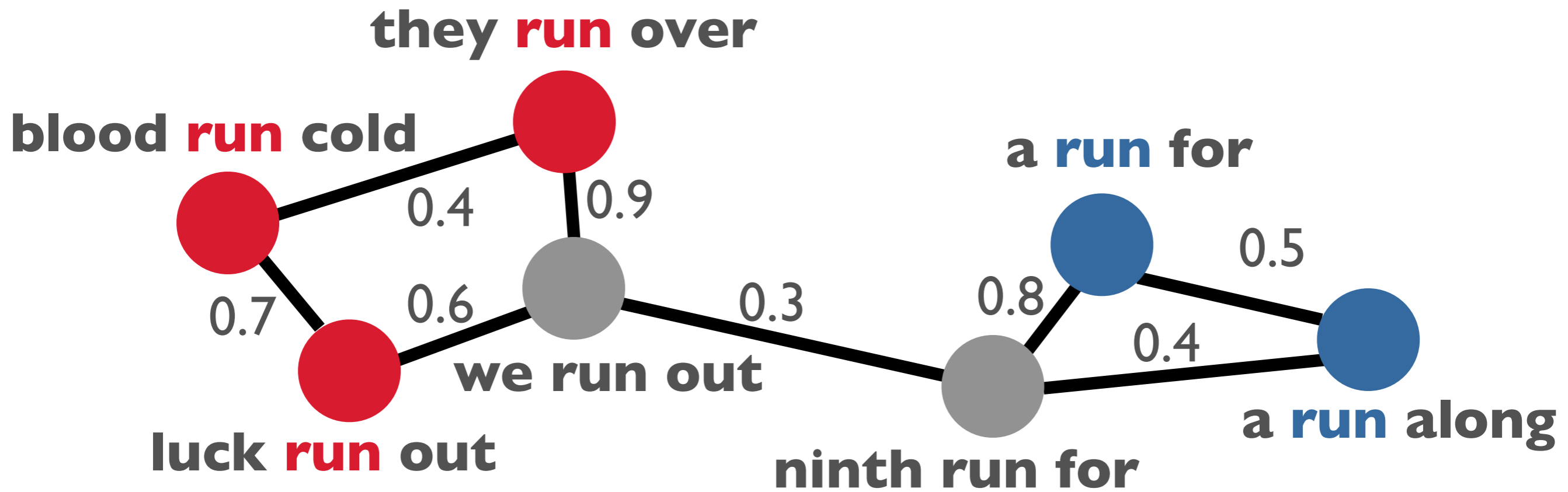
GRAPH-BASED LEARNING

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



GRAPH-BASED LEARNING

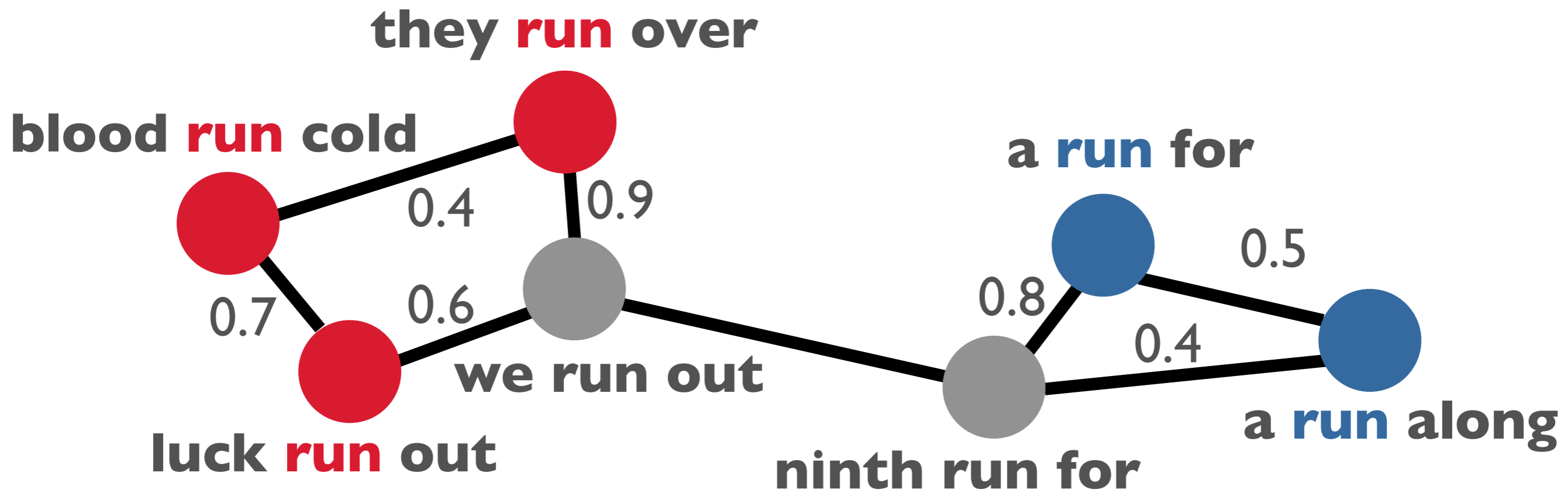
Labels: **verb (V)**, **noun (N)**, etc.



$$\| \begin{matrix} \text{N} \\ \text{V} \end{matrix} - \begin{matrix} \text{V} \\ \text{N} \end{matrix} \|_2^2$$

GRAPH-BASED LEARNING

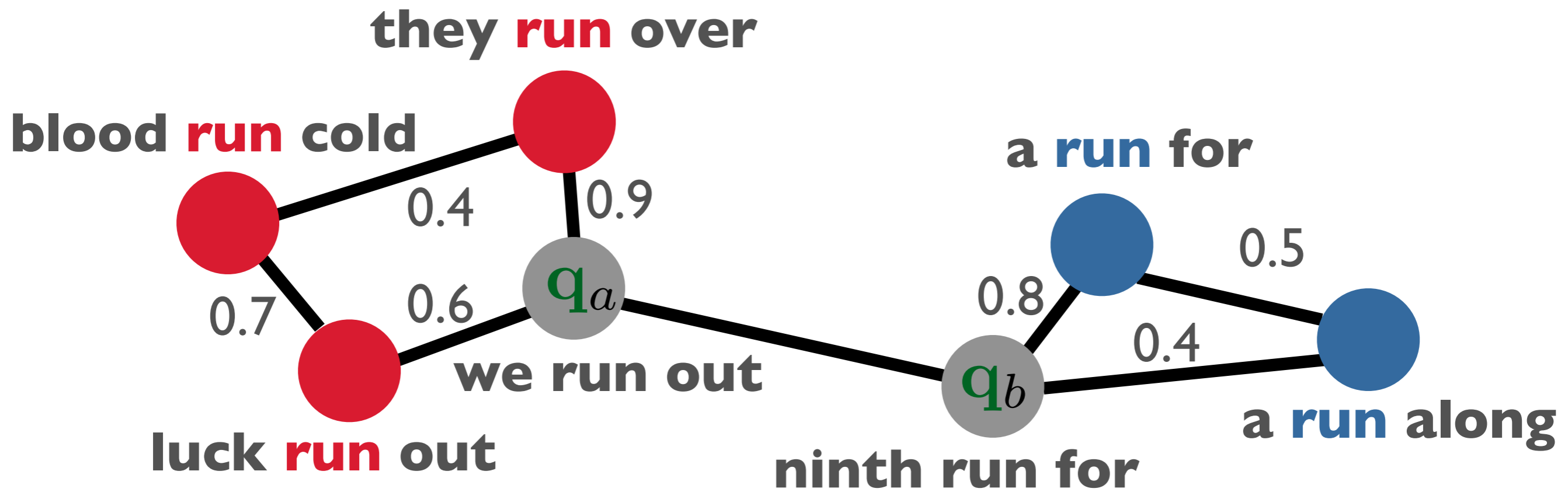
Labels: **verb (V)**, **noun (N)**, etc.



$$0.3 \quad || \quad \begin{array}{c} \text{N} \\ \text{V} \end{array} \quad - \quad \begin{array}{c} \text{V} \\ \text{N} \end{array} \quad || \quad 2$$

GRAPH-BASED LEARNING

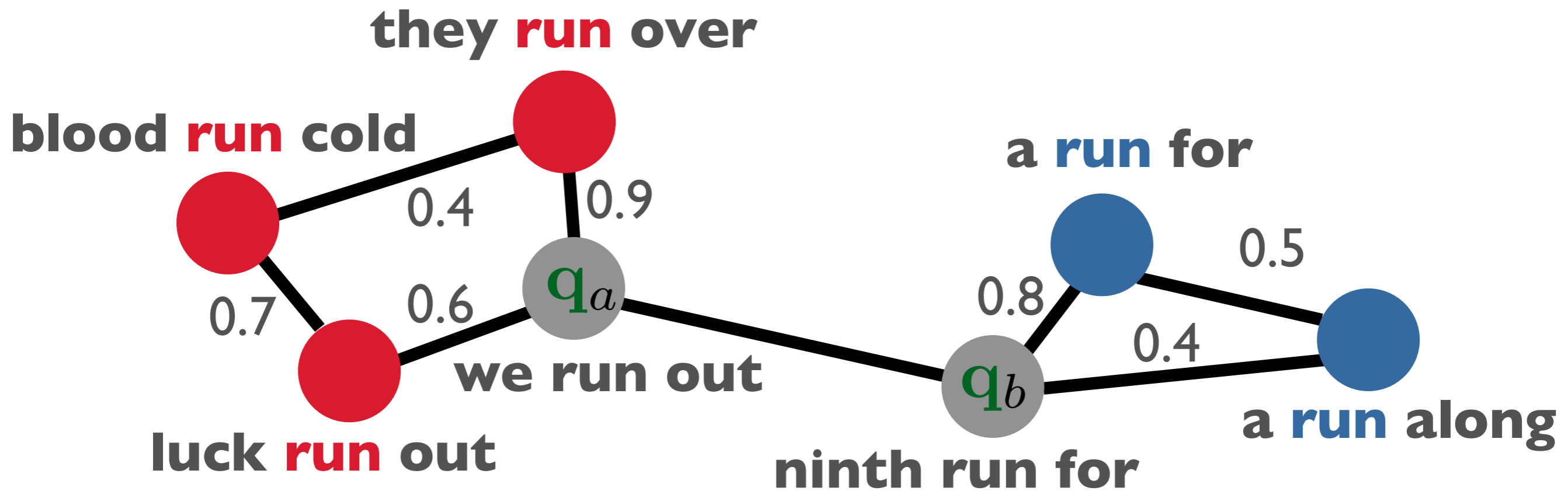
Labels: **verb (V)**, **noun (N)**, etc.



$$w_{ab} \propto \|q_a - q_b\|_2^2$$

GRAPH-BASED LEARNING

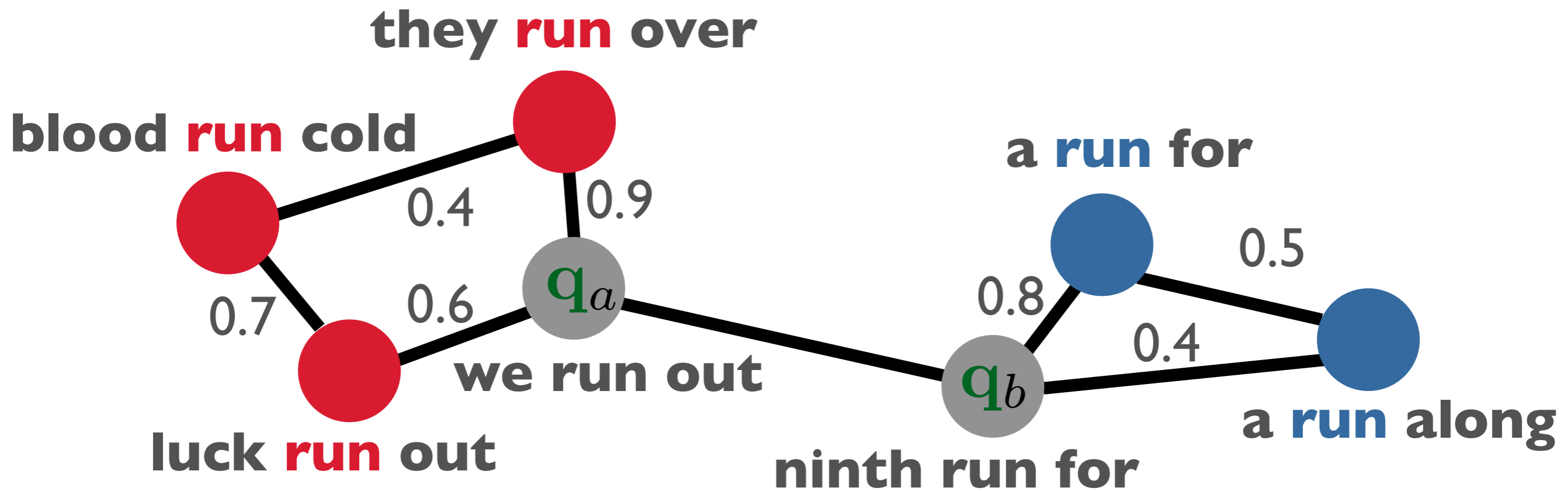
Labels: **verb (V)**, **noun (N)**, etc.



$$\text{Lap}(q) = w_{ab} \left\| \mathbf{q}_a - \mathbf{q}_b \right\|_2^2$$

GRAPH-BASED LEARNING

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



$$\text{Lap}(q) = \sum_{a=1}^N \sum_{b=L+1}^N w_{ab} \|\mathbf{q}_a - \mathbf{q}_b\|_2^2$$

STRUCTURED PREDICTION

ninth run for

STRUCTURED PREDICTION

ninth run for

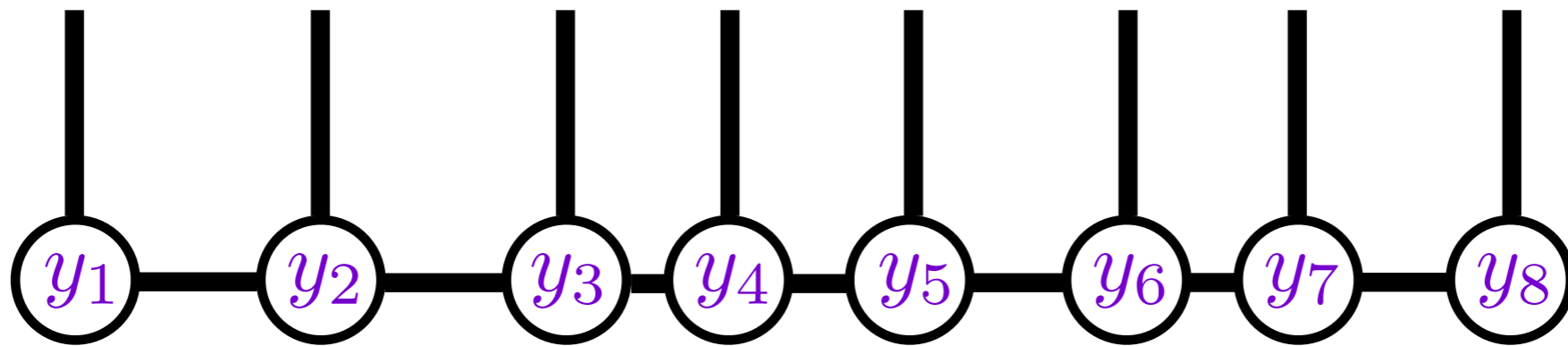
STRUCTURED PREDICTION

The soldiers of the ninth run for cover

STRUCTURED PREDICTION

The soldiers of the ninth run for cover

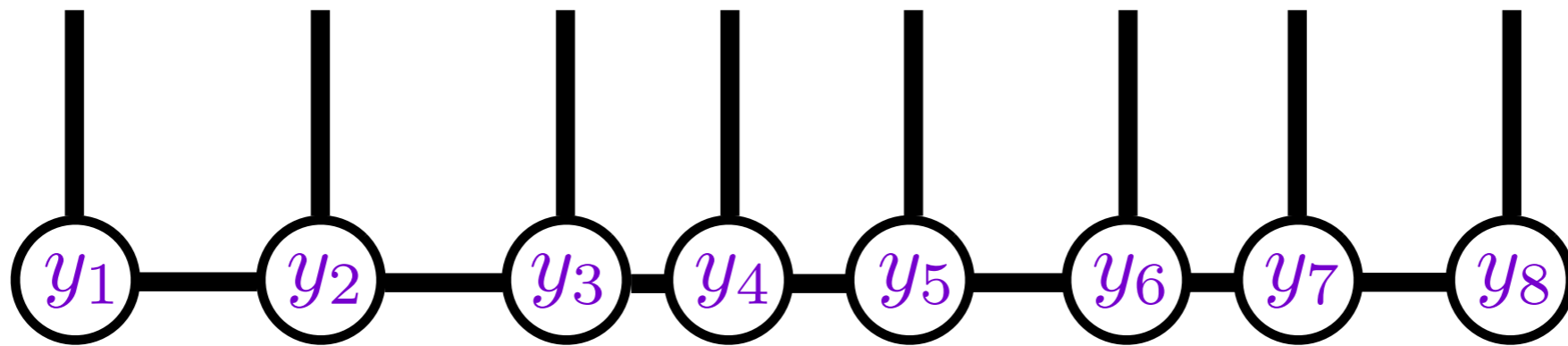
CRF



STRUCTURED PREDICTION

The soldiers of the ninth run for cover

CRF

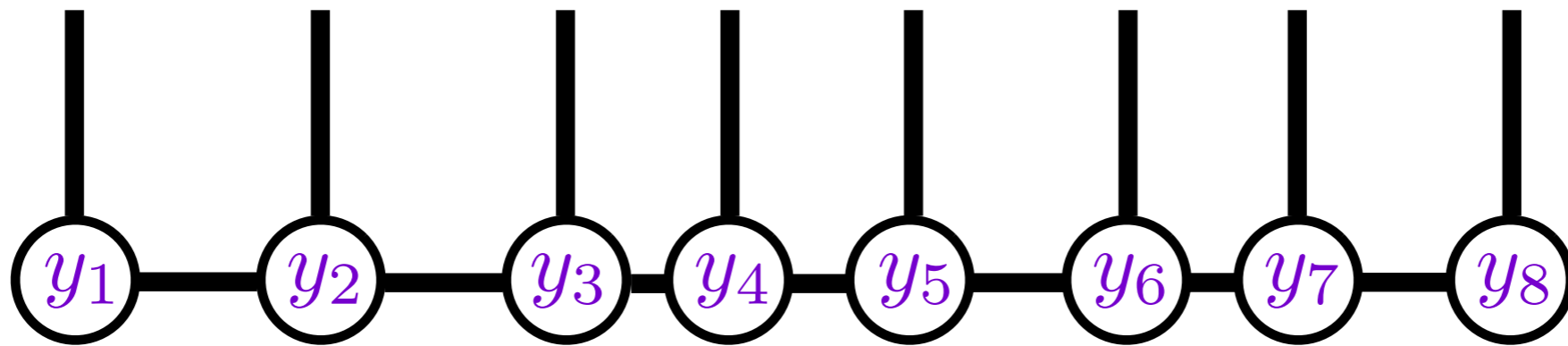


$f(\quad)$

STRUCTURED PREDICTION

The soldiers of the ninth run for cover

CRF

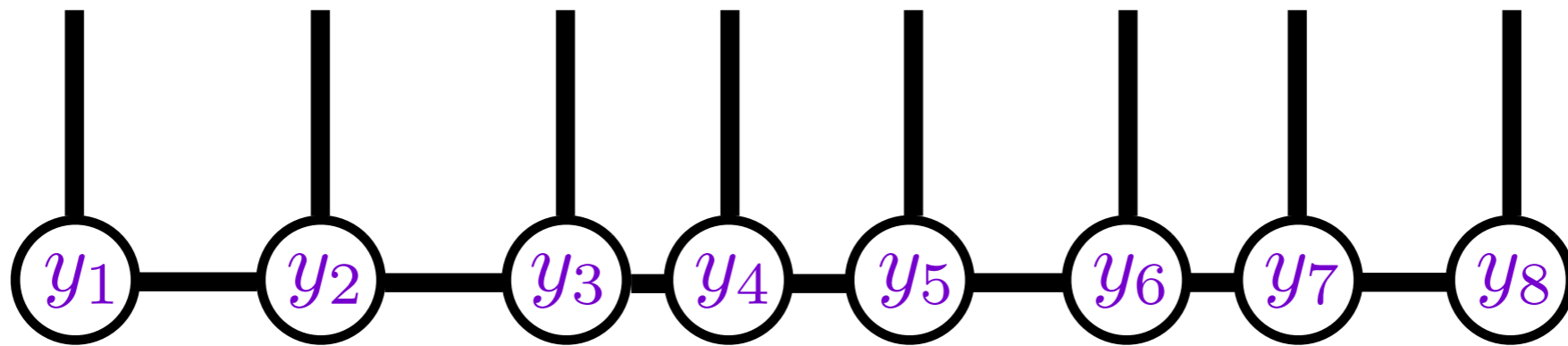


$$\mathbf{f}(y_t)$$

STRUCTURED PREDICTION

The soldiers of the ninth run for cover

CRF

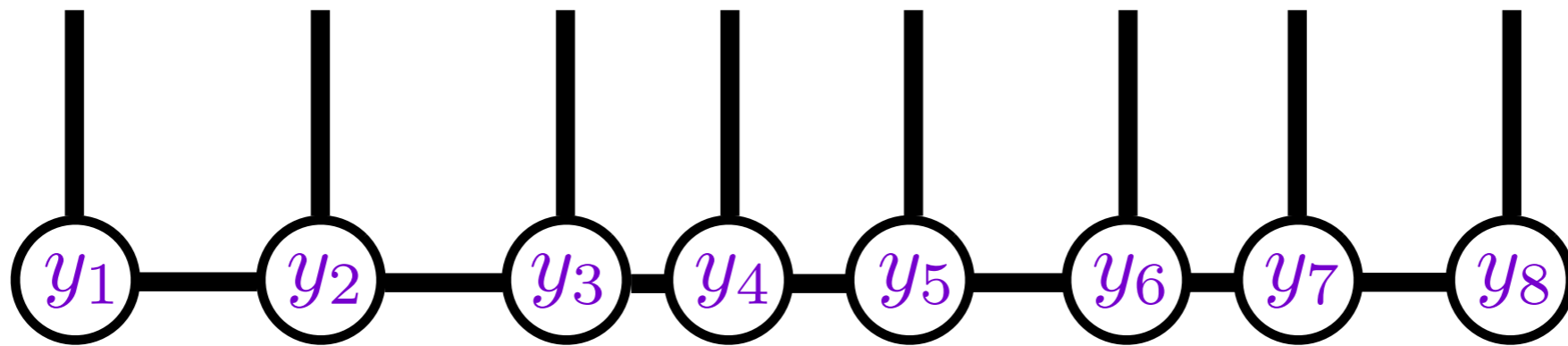


$$\mathbf{f}(y_t, y_{t-1})$$

STRUCTURED PREDICTION

$\mathbf{x} =$ **The soldiers of the ninth run for cover**

CRF

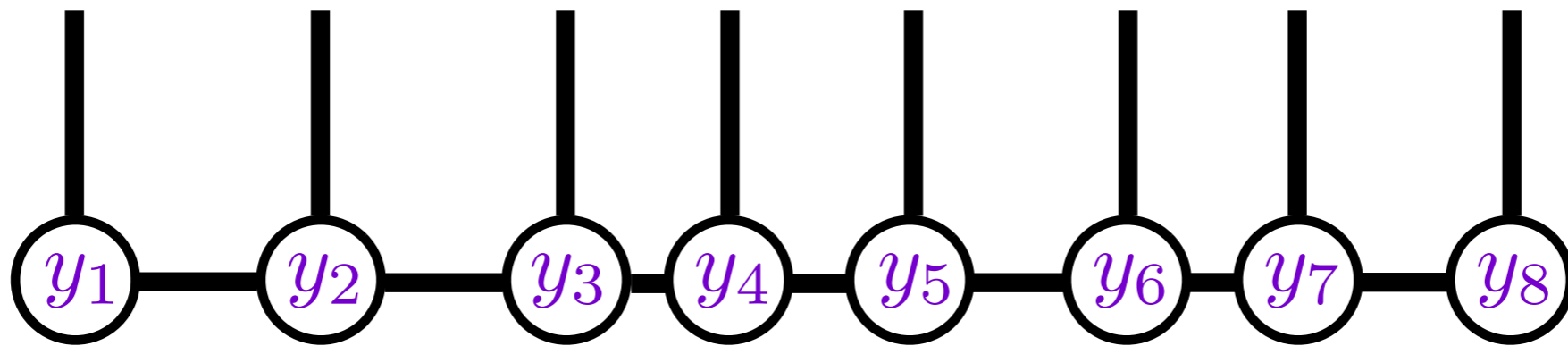


$$\mathbf{f}(y_t, y_{t-1}, \mathbf{x})$$

STRUCTURED PREDICTION

$\mathbf{x} =$ **The soldiers of the ninth run for cover**

CRF



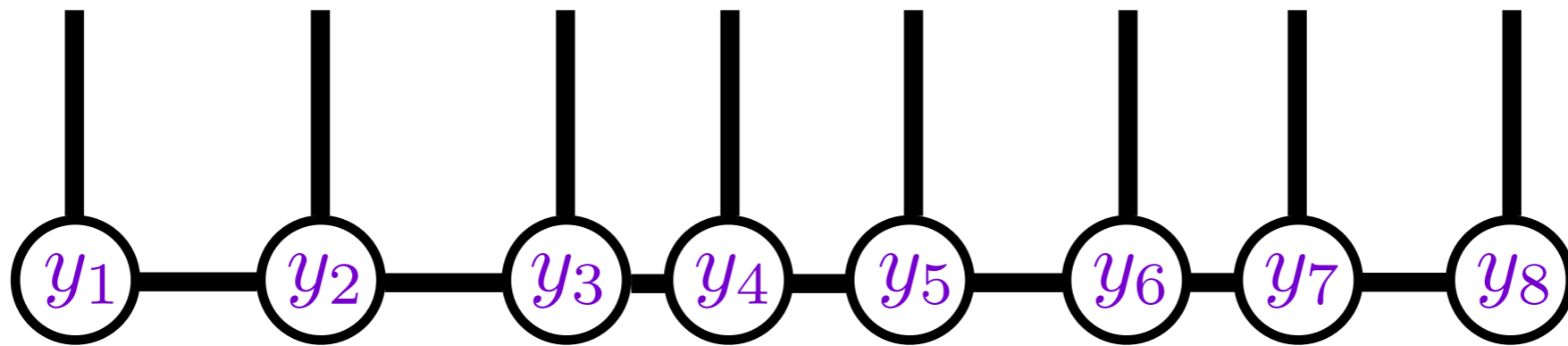
$$\mathbf{f}(y_t, y_{t-1}, \mathbf{x})$$

$\underbrace{\hspace{10em}}$
 p -factor

STRUCTURED PREDICTION

$\mathbf{x} =$ **The soldiers of the ninth run for cover**

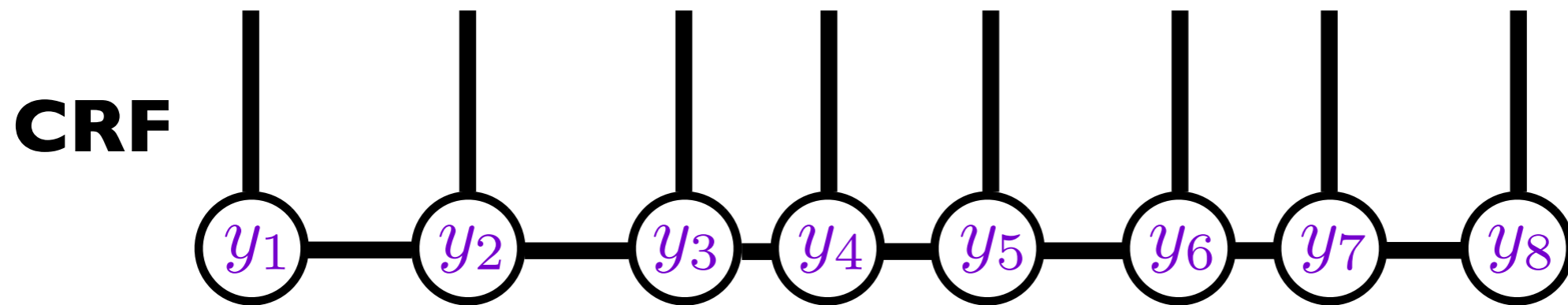
CRF



$$p_{\theta}(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z_{\theta}(\mathbf{x})} \exp \left[\sum_{t=1}^T \theta^{\top} \underbrace{\mathbf{f}(y_t, y_{t-1}, \mathbf{x})}_{p\text{-factor}} \right]$$

STRUCTURED PREDICTION

$\mathbf{x} =$ **The soldiers of the ninth run for cover**



$$p_{\theta}(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z_{\theta}(\mathbf{x})} \exp \left[\sum_{t=1}^T \theta^{\top} \underbrace{\mathbf{f}(y_t, y_{t-1}, \mathbf{x})}_{p\text{-factor}} \right]$$

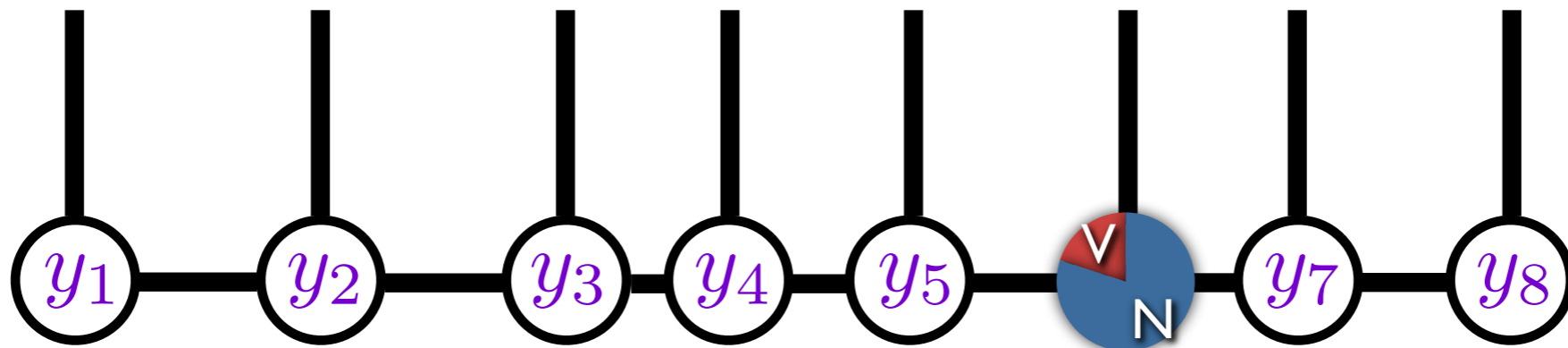
$$\text{NLik}(p_{\theta}) = - \sum_{i=1}^{\ell} \log p_{\theta}(\mathbf{y}^i \mid \mathbf{x}^i)$$

STRUCTURED PREDICTION

$\mathbf{x} =$ **The soldiers of the**

cover

CRF



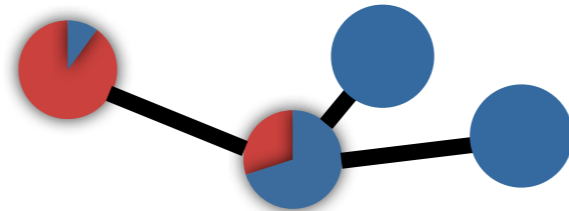
$$p_{\theta}(y_t | \mathbf{x})$$

WHY COMBINE?

Each type of learning incorporates different information

WHY COMBINE?

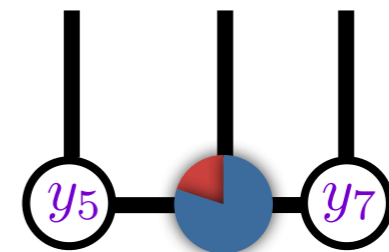
Each type of learning incorporates different information



ninth run for

graph-propagation

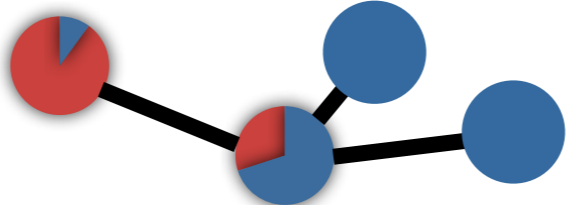
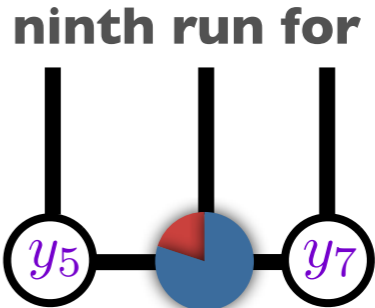
ninth run for



CRF estimation

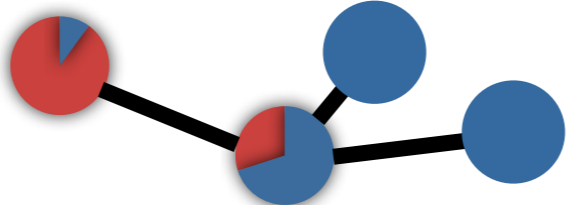
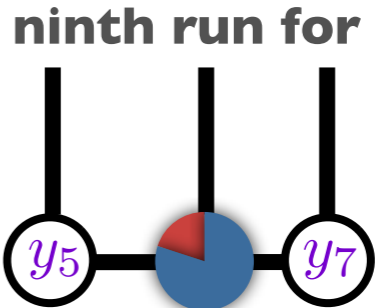
WHY COMBINE?

Each type of learning incorporates different information

	 <p>ninth run for graph-propagation</p>	 <p>ninth run for CRF estimation</p>
Data	unlabeled	labeled

WHY COMBINE?

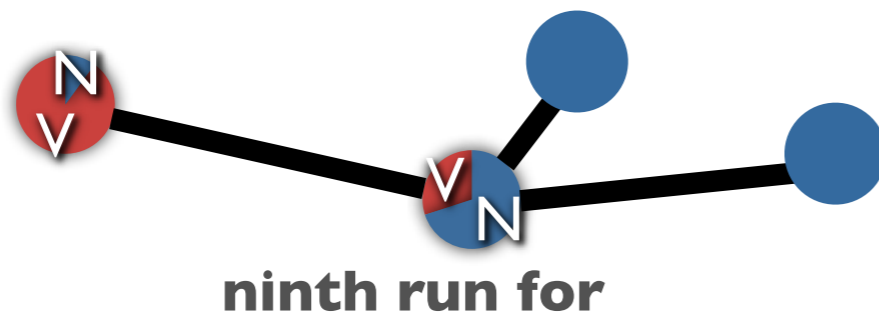
Each type of learning incorporates different information

	 <p>ninth run for graph-propagation</p>	 <p>ninth run for CRF estimation</p>
Data	unlabeled	labeled
Context	trigram	sentence

PRIOR WORK

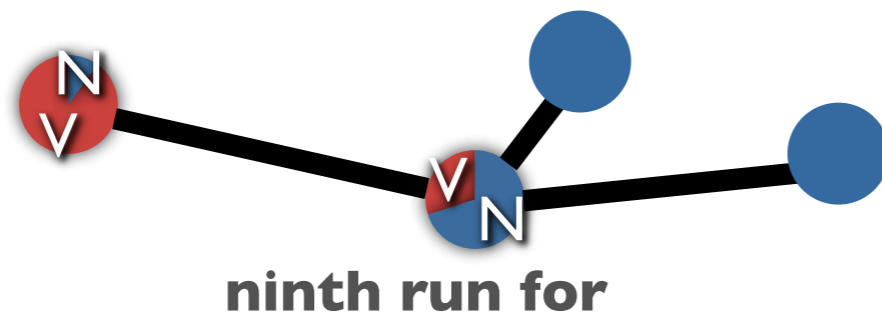
PRIOR WORK

Lap(q) graph-propagation

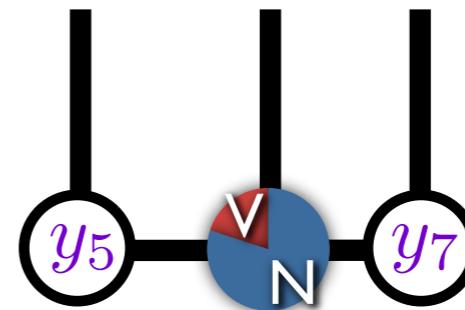


PRIOR WORK

$\text{Lap}(q)$ graph-propagation + CRF estimation $\text{NLik}(p_\theta)$



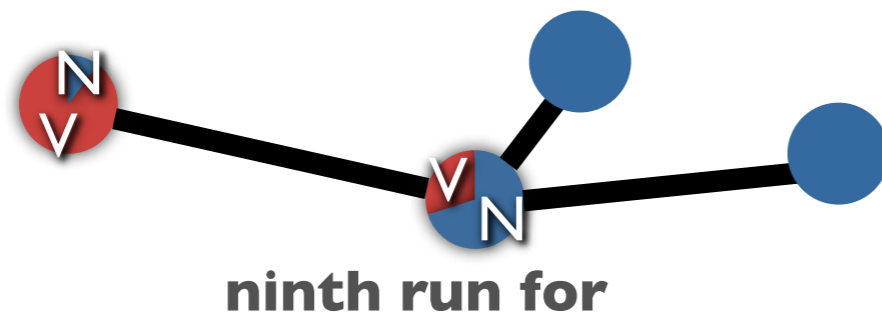
ninth run for



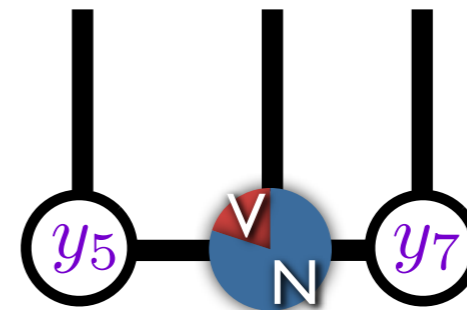
PRIOR WORK

Subramanya et al. (EMNLP 2010)

$\text{Lap}(q)$ graph-propagation + CRF estimation $\text{NLik}(p_\theta)$



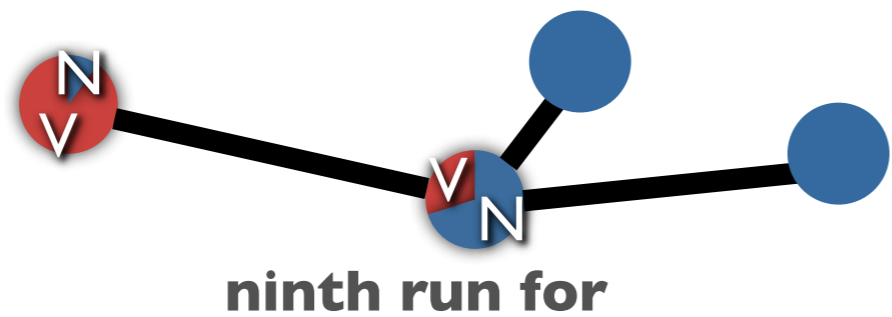
ninth run for



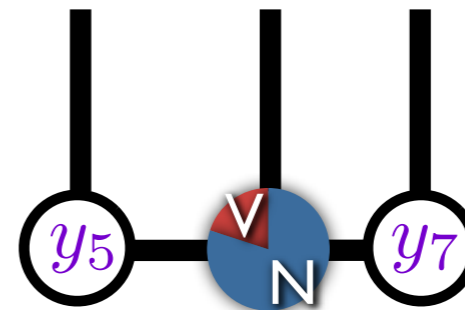
PRIOR WORK

Subramanya et al. (EMNLP 2010)

$\text{Lap}(q)$ graph-propagation + CRF estimation $\text{NLik}(p_\theta)$



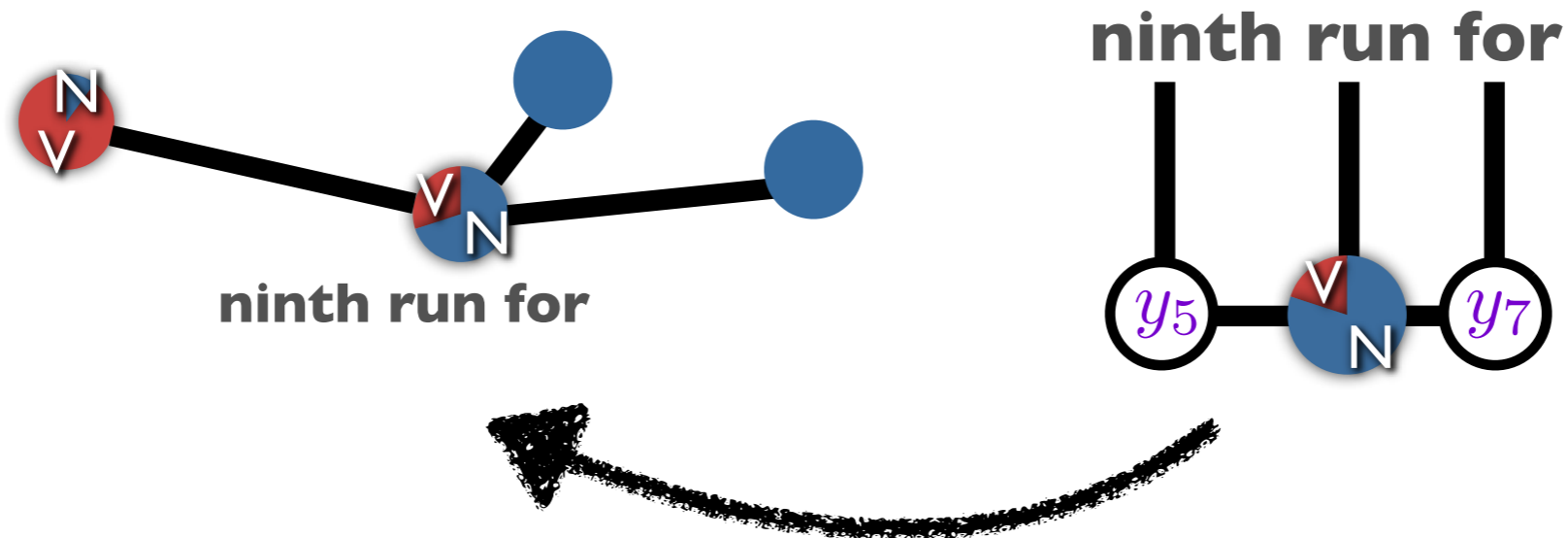
ninth run for



PRIOR WORK

Subramanya et al. (EMNLP 2010)

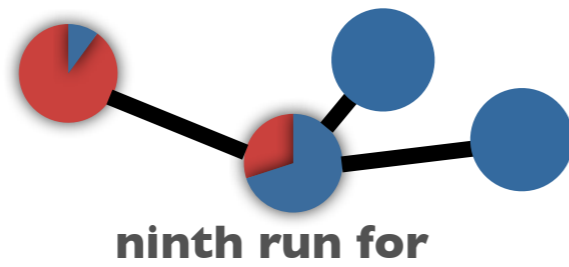
$\text{Lap}(q)$ graph-propagation + CRF estimation $\text{NLik}(p_\theta)$



This work: retains efficiency while optimizing an extendible, joint objective.

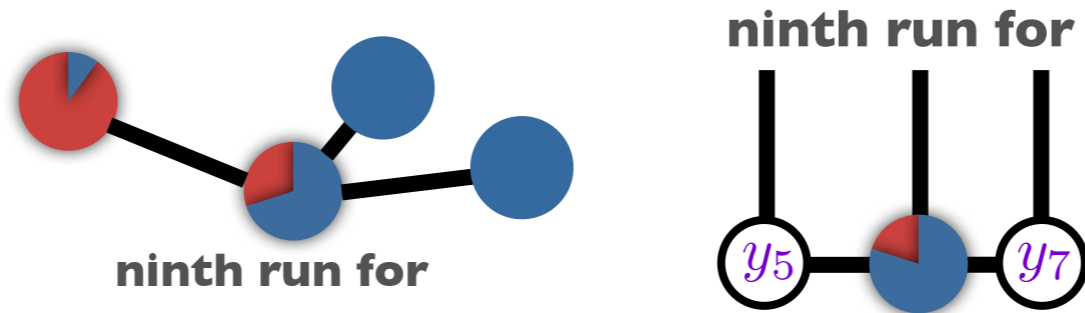
JOINT OBJECTIVE

JOINT OBJECTIVE



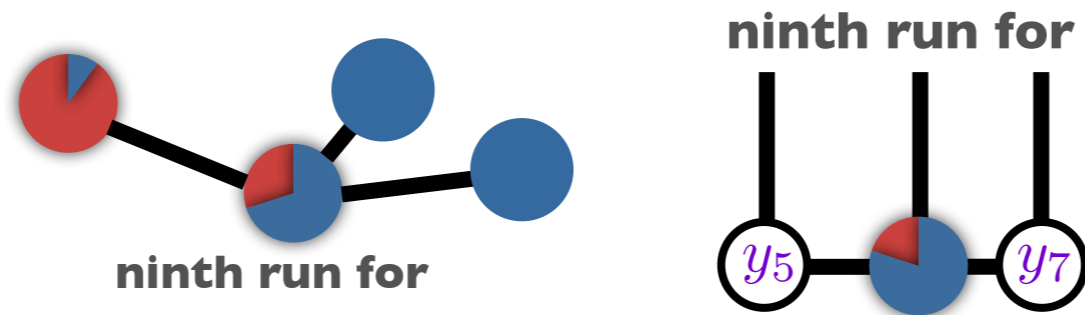
Lap(q)

JOINT OBJECTIVE



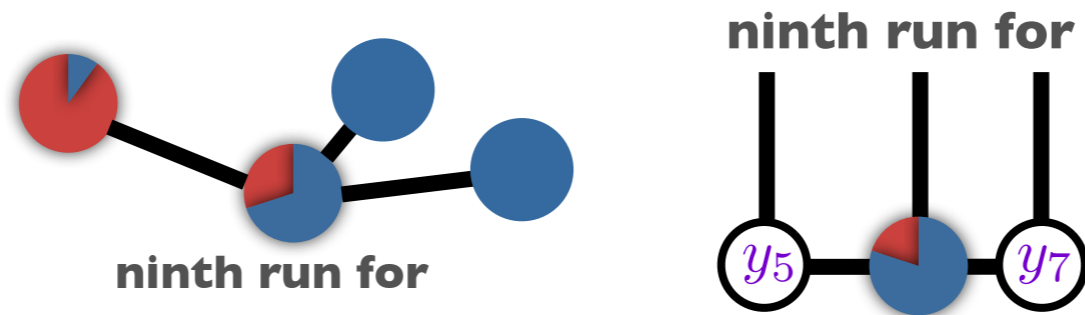
$$\text{Lap}(q) + \text{NLik}(p_\theta)$$

JOINT OBJECTIVE



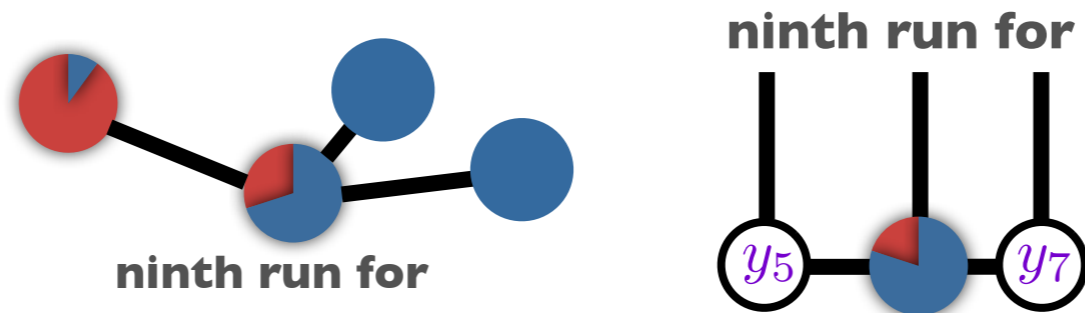
$$\mathcal{J}(q, p_\theta) = \text{Lap}(q) + \text{NLik}(p_\theta)$$

JOINT OBJECTIVE



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

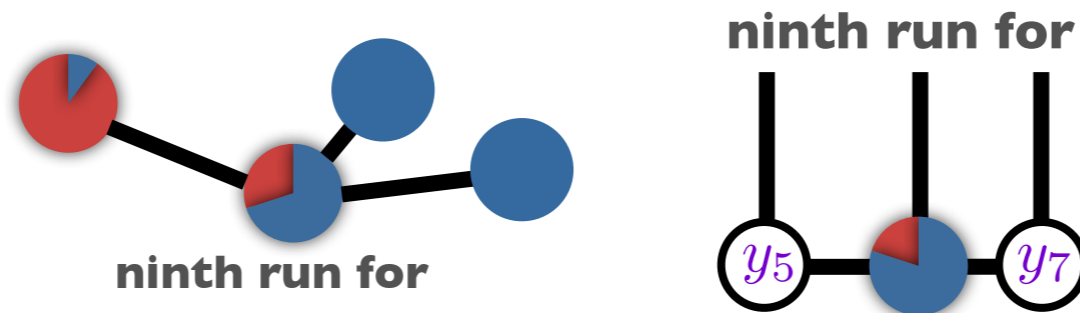
JOINT OBJECTIVE



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

The soldiers of the ninth run for cover

JOINT OBJECTIVE

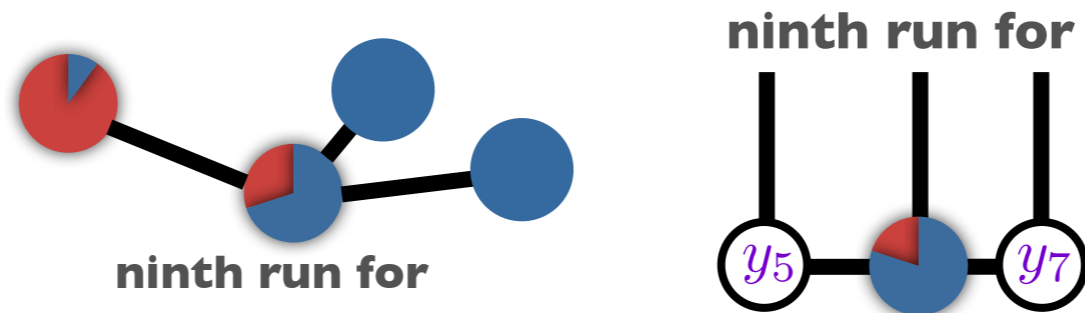


$$\mathcal{J}(q, p_\theta) = \text{Lap}(q) + \text{NLik}(p_\theta) + \text{KL}(q \parallel p_\theta)$$

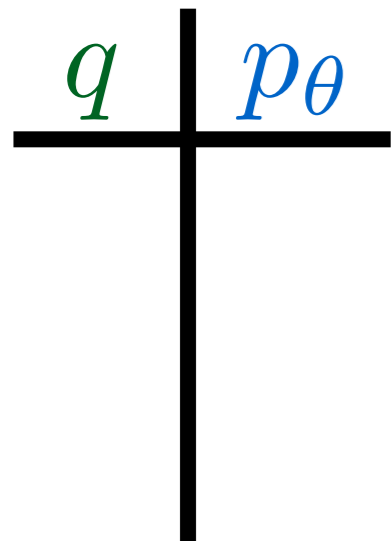
The soldiers of the ninth run for cover

$(\# \text{ tags})^8 \left\{ \begin{array}{cccccccccc} \mathbf{N} & & \mathbf{N} & & \mathbf{N} & & \mathbf{N} & & \mathbf{N} & & \mathbf{N} \\ \mathbf{N} & & \mathbf{N} & & \mathbf{N} & & \mathbf{N} & & \mathbf{N} & & \mathbf{V} \\ & & & & \dots & & & & & & \end{array} \right.$

JOINT OBJECTIVE



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

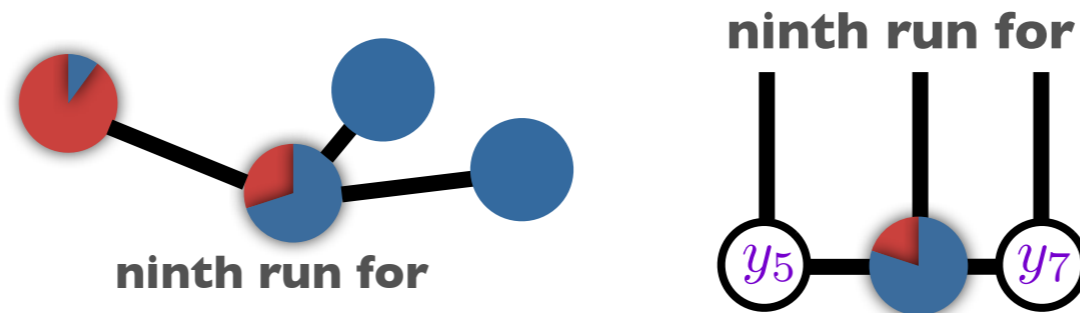


The soldiers of the ninth run for cover

N N N N N N N N N
N N N N N N N N V

...

JOINT OBJECTIVE



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

q	p_θ
7e-5	2e-5
3e-6	8e-6
...	...

The soldiers of the ninth run for cover

N	N	N	N	N	N	N	N	N
N	N	N	N	N	N	N	N	V
			...					

OPTIMIZATION

$$\min_{q, \theta} \mathcal{J}(q, p_{\theta})$$


OPTIMIZATION

$$\min_{q, \theta} \mathcal{J}(q, p_{\theta})$$

← unconstrained

OPTIMIZATION

$$\min_{q, \theta} \mathcal{J}(q, p_{\theta})$$

Δ  unconstrained

OPTIMIZATION

$$\min_{q, \theta} \mathcal{J}(q, p_{\theta})$$

Δ \swarrow \nwarrow unconstrained

p update:

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

OPTIMIZATION

$$\min_{q, \theta} \mathcal{J}(q, p_{\theta})$$

Δ \swarrow \nwarrow unconstrained

p update:

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

Next 3 slides: Why several common techniques don't work for updating q

OPTIMIZATION

$$\min_{q, \theta} \mathcal{J}(q, p_{\theta})$$

OPTIMIZATION

$$\min_{q, \theta} \mathcal{J}(q, p_{\theta})$$

q update:

$$q_{\mathbf{y}}^{i'} = q_{\mathbf{y}}^i - \eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial q_{\mathbf{y}}^i}$$

OPTIMIZATION

$$\min_{q, \theta} \mathcal{J}(q, p_{\theta})$$

q update:

$$q_{\mathbf{y}}^{i'} = \text{proj}_{\Delta} \left(q_{\mathbf{y}}^i - \eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial q_{\mathbf{y}}^i} \right)$$

OPTIMIZATION

$$\min_{q, \theta} \mathcal{J}(q, p_{\theta})$$

q update:

$$q_{\mathbf{y}}^{i'} = \text{proj}_{\Delta} \left(q_{\mathbf{y}}^i - \eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial q_{\mathbf{y}}^i} \right)$$

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

OPTIMIZATION

$$\min_{q, \theta} \mathcal{J}(q, p_{\theta})$$

q update:

$$q_{\mathbf{y}}^{i'} = \text{proj}_{\Delta} \left(q_{\mathbf{y}}^i - \eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial q_{\mathbf{y}}^i} \right)$$

$q^i \in \Delta$ of dimension $(\# \text{ tags})^{(i\text{'s length})}$

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

OPTIMIZATION

$$\min_{q, \theta} \mathcal{J}(q, p_{\theta})$$

q update:

$$q_{\mathbf{y}}^{i'} = \text{proj}_{\Delta} \left(q_{\mathbf{y}}^i - \eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial q_{\mathbf{y}}^i} \right)$$

$q^i \in \Delta$ of dimension $(\# \text{ tags})^{(i\text{'s length})}$

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

-Problem 2: no compact form $(\# \text{ tags})^{(i\text{'s length})}$ values

OPTIMIZATION

$$\min_{q, \theta} \mathcal{J}(q, p_{\theta})$$

q update:

~~$$q_{\mathbf{y}}^{i'} = \text{proj}_{\Delta} \left(q_{\mathbf{y}}^i - \eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial q_{\mathbf{y}}^i} \right)$$~~

$q^i \in \Delta$ of dimension $(\# \text{ tags})^{(i\text{'s length})}$

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

-Problem 2: no compact form $(\# \text{ tags})^{(i\text{'s length})}$ values

DUAL OPTIMIZATION

$$\mathcal{J}(q, p_\theta)$$

DUAL OPTIMIZATION

$$\mathcal{J}(q, p_{\theta}) + \gamma \left(\sum_{\mathbf{y}} q_{\mathbf{y}}^i - 1 \right)$$

DUAL OPTIMIZATION

$$\mathcal{J}(q, p_{\theta}) + \gamma \left(\sum_{\mathbf{y}} q_{\mathbf{y}}^i - 1 \right)$$

Posterior Regularization (PR) uses dual
Ganchev et al. (JMLR 2010)

DUAL OPTIMIZATION

$$\mathcal{J}(q, p_{\theta}) + \gamma \left(\sum_{\mathbf{y}} q_{\mathbf{y}}^i - 1 \right)$$

Posterior Regularization (PR) uses dual
Ganchev et al. (JMLR 2010)

This work: $\text{Lap}(q)$

DUAL OPTIMIZATION

$$\mathcal{J}(q, p_{\theta}) + \gamma \left(\sum_{\mathbf{y}} q_{\mathbf{y}}^i - 1 \right)$$

Posterior Regularization (PR) uses dual
Ganchev et al. (JMLR 2010)

This work: $\text{Lap}(q)$ \longrightarrow Standard PR: simpler

DUAL OPTIMIZATION

$$\mathcal{J}(q, p_{\theta}) + \gamma \left(\sum_{\mathbf{y}} q_{\mathbf{y}}^i - 1 \right)$$

Posterior Regularization (PR) uses dual
Ganchev et al. (JMLR 2010)

This work: $\text{Lap}(q)$ \longrightarrow Standard PR: simpler
 \uparrow
 p -factors
 y_t, y_{t-1}, \mathbf{x}

DUAL OPTIMIZATION

$$\mathcal{J}(q, p_{\theta}) + \gamma \left(\sum_{\mathbf{y}} q_{\mathbf{y}}^i - 1 \right)$$

Posterior Regularization (PR) uses dual
Ganchev et al. (JMLR 2010)

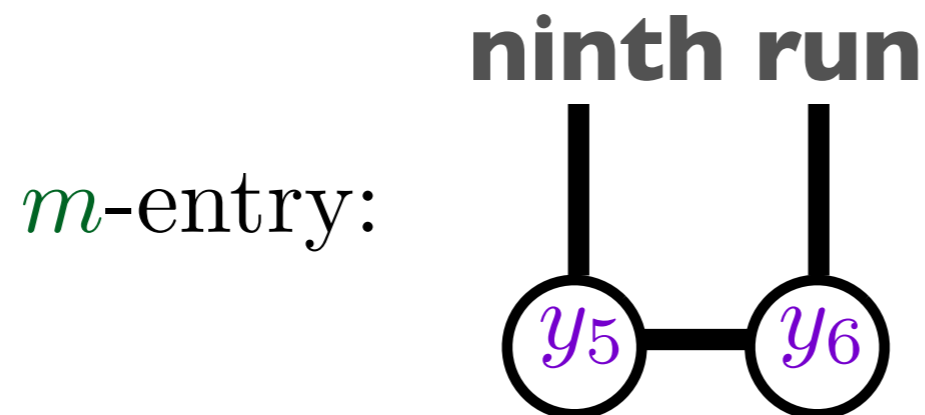
This work: **Lap**(q) \longrightarrow Standard PR: **Linear**(m)

DUAL OPTIMIZATION

$$\mathcal{J}(q, p_{\theta}) + \gamma \left(\sum_{\mathbf{y}} q_{\mathbf{y}}^i - 1 \right)$$

Posterior Regularization (PR) uses dual
Ganchev et al. (JMLR 2010)

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



DUAL OPTIMIZATION

$$\mathcal{J}(q, p_{\theta}) + \gamma \left(\sum_{\mathbf{y}} q_{\mathbf{y}}^i - 1 \right)$$

Posterior Regularization (PR) uses dual
Ganchev et al. (JMLR 2010)

This work: **Lap**(q) \longrightarrow Standard PR: **Linear**(m)

DUAL OPTIMIZATION

$$\mathcal{J}(q, p_{\theta}) + \gamma \left(\sum_{\mathbf{y}} q_{\mathbf{y}}^i - 1 \right)$$

Posterior Regularization (PR) uses dual
Ganchev et al. (JMLR 2010)

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

\swarrow
 $\text{Lap}(m)$

DUAL OPTIMIZATION

$$\mathcal{J}(q, p_{\theta}) + \gamma \left(\sum_{\mathbf{y}} q_{\mathbf{y}}^i - 1 \right)$$

Posterior Regularization (PR) uses dual
Ganchev et al. (JMLR 2010)

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

\swarrow
 $\text{Lap}(m)$, a quadratic function

DUAL OPTIMIZATION

$$\mathcal{J}(q, p_\theta) + \gamma \left(\sum_{\mathbf{y}} q_{\mathbf{y}}^i - 1 \right)$$

Posterior Regularization (PR) uses dual
Ganchev et al. (JMLR 2010)

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

$\text{Lap}(m)$, a quadratic function
Dual of quadratic requires:

$$\begin{pmatrix} 1 & 2 & \dots & N \\ \begin{matrix} 1 \\ 2 \\ \vdots \\ N \end{matrix} & \text{Red Square} \end{pmatrix}^{-1}$$

~~DUAL OPTIMIZATION~~

$$\mathcal{J}(q, p_{\theta}) + \gamma \left(\sum_{\mathbf{y}} q_{\mathbf{y}}^i - 1 \right)$$

Posterior Regularization (PR) uses dual
Ganchev et al. (JMLR 2010)

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

$\text{Lap}(m)$, a quadratic function
Dual of quadratic requires:

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

EXPONENTIATED GRADIENT

EXPONENTIATED GRADIENT

$$q_{\mathbf{y}}^{i'} \propto q_{\mathbf{y}}^i \exp \left[-\eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial q_{\mathbf{y}}^i} \right]$$

EXPONENTIATED GRADIENT

$$q_{\mathbf{y}}^{i'} \propto q_{\mathbf{y}}^i \exp \left[-\eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial q_{\mathbf{y}}^i} \right]$$

Collins et al. (JMLR 2008): Exponentiated gradient for CRFs

EXPONENTIATED GRADIENT

$$q_{\mathbf{y}}^{i'} \propto q_{\mathbf{y}}^i \exp \left[-\eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial q_{\mathbf{y}}^i} \right]$$

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

EXPONENTIATED GRADIENT

$$q_{\mathbf{y}}^{i'} \propto q_{\mathbf{y}}^i \exp \left[-\eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial q_{\mathbf{y}}^i} \right]$$

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

$$\exp \left[-\eta \sum_{t=1}^T \frac{\partial \text{Lap}(m_{\mathbf{y}}^i)}{\partial m_{t, y_t, y_{t-1}}^i} \right]$$

EXPONENTIATED GRADIENT

$$q_{\mathbf{y}}^{i'} \propto q_{\mathbf{y}}^i \exp \left[-\eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial q_{\mathbf{y}}^i} \right]$$

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

$$\exp \left[-\eta \sum_{t=1}^T \frac{\partial \text{Lap}(m_{\mathbf{y}}^i)}{\partial m_{t, y_t, y_{t-1}}^i} \right]$$

product of p-factors

EXPONENTIATED GRADIENT

$$q_{\mathbf{y}}^{i'} \propto q_{\mathbf{y}}^i \exp \left[-\eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial q_{\mathbf{y}}^i} \right]$$

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

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

product of p-factors

EXPONENTIATED GRADIENT

$$q_{\mathbf{y}}^{i'} \propto q_{\mathbf{y}}^i \exp \left[-\eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial q_{\mathbf{y}}^i} \right]$$

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

$$\exp \left[-\eta \sum_{t=1}^T \frac{\partial \text{Lap}(m_{\mathbf{y}}^i)}{\partial m_{t, y_t, y_{t-1}}^i} \right] \underbrace{p_{\theta}(\mathbf{y} \mid \mathbf{x}^i)^{\eta} (q_{\mathbf{y}}^i)^{-\eta} e}_{\text{product of p-factors}}$$

EXPONENTIATED GRADIENT

$$q_{\mathbf{y}}^{i'} \propto q_{\mathbf{y}}^i \exp \left[-\eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial q_{\mathbf{y}}^i} \right]$$

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

$$\exp \left[-\eta \sum_{t=1}^T \frac{\partial \text{Lap}(m_{\mathbf{y}}^i)}{\partial m_{t, y_t, y_{t-1}}^i} \right] \underbrace{p_{\theta}(\mathbf{y} \mid \mathbf{x}^i)^{\eta} (q_{\mathbf{y}}^i)^{-\eta} e}_{\text{product of p-factors}}$$

product of p-factors

EXPONENTIATED GRADIENT

$$q_{\mathbf{y}}^{i'} \propto q_{\mathbf{y}}^i \exp \left[-\eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial q_{\mathbf{y}}^i} \right]$$

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

$$\exp \left[-\eta \sum_{t=1}^T \frac{\partial \text{Lap}(m_{\mathbf{y}}^i)}{\partial m_{t, y_t, y_{t-1}}^i} \right] \underbrace{p_{\theta}(\mathbf{y} \mid \mathbf{x}^i)^{\eta} (q_{\mathbf{y}}^i)^{-\eta} e}_{\text{product of p-factors}}$$

proj_{Δ}

EXPONENTIATED GRADIENT

$$q_{\mathbf{y}}^{i'} \propto q_{\mathbf{y}}^i \exp \left[-\eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial q_{\mathbf{y}}^i} \right]$$

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

$$\exp \left[-\eta \sum_{t=1}^T \frac{\partial \text{Lap}(m_{\mathbf{y}}^i)}{\partial m_{t, y_t, y_{t-1}}^i} \right] \underbrace{p_{\theta}(\mathbf{y} \mid \mathbf{x}^i)^{\eta} (q_{\mathbf{y}}^i)^{-\eta} e}_{\text{product of p-factors}}$$

$$\text{proj}_{\Delta} \longrightarrow Z_q(\mathbf{x}^i)$$

EXPONENTIATED GRADIENT

$$q_{\mathbf{y}}^{i'} \propto q_{\mathbf{y}}^i \exp \left[-\eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial q_{\mathbf{y}}^i} \right]$$

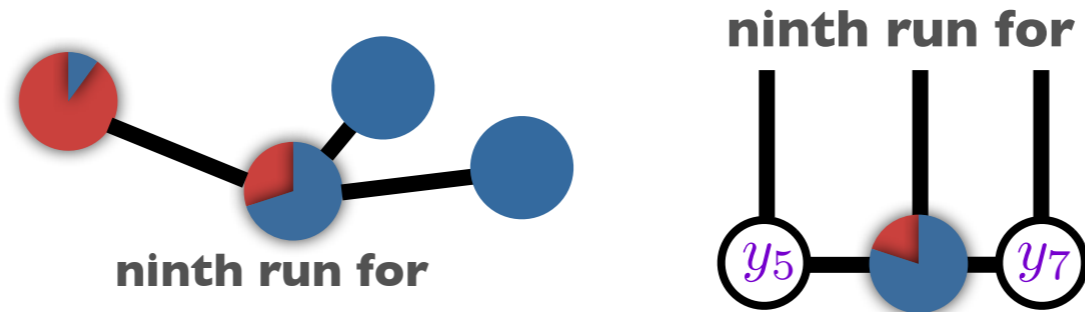
$$\exp \left[-\eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial q_{\mathbf{y}}^i} \right] =$$

$$\exp \left[-\eta \sum_{t=1}^T \frac{\partial \text{Lap}(m_{\mathbf{y}}^i)}{\partial m_{t, y_t, y_{t-1}}^i} \right] \underbrace{p_{\theta}(\mathbf{y} \mid \mathbf{x}^i)^{\eta} (q_{\mathbf{y}}^i)^{-\eta} e}_{\text{product of p-factors}}$$

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

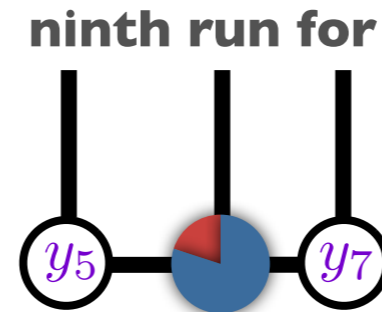
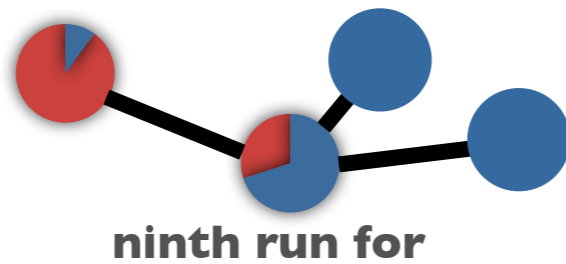
SUMMARY

SUMMARY



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

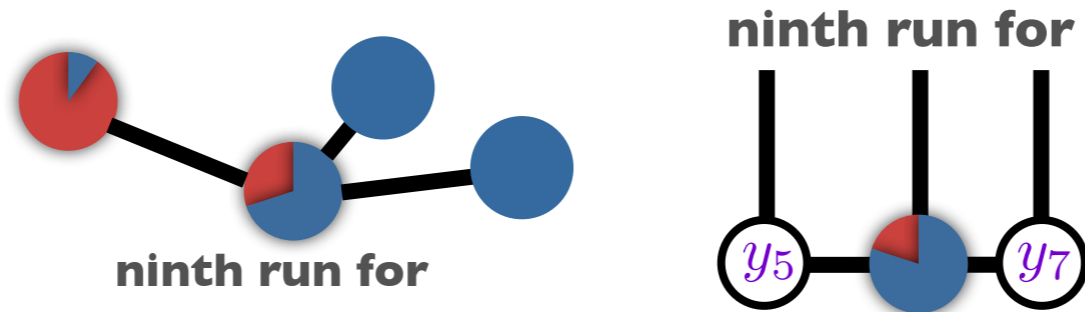
SUMMARY



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

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

SUMMARY

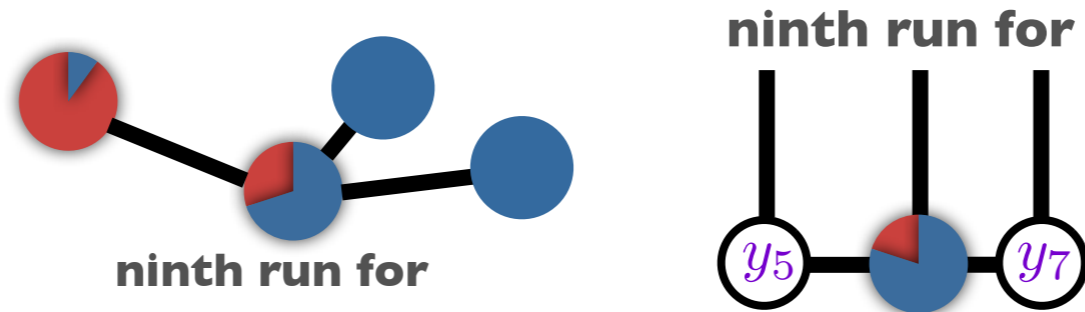


$$\mathcal{J}(q, p_\theta) = \text{Lap}(q) + \text{NLik}(p_\theta) + \text{KL}(q \parallel p_\theta)$$

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

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

SUMMARY

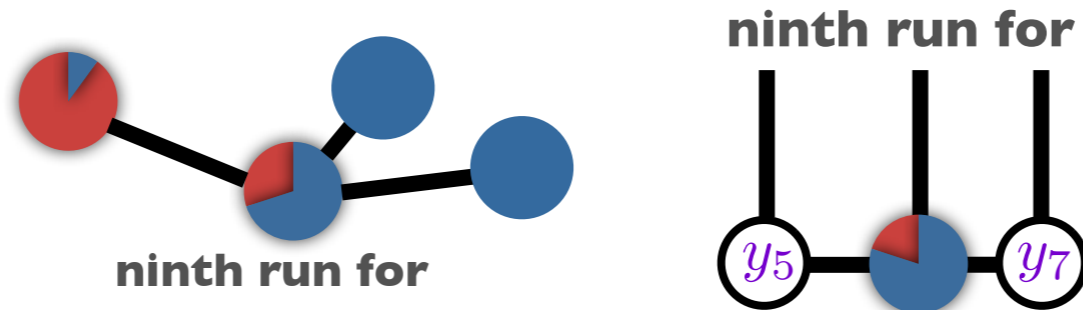


$$\mathcal{J}(q, p_\theta) = \text{Lap}(q) + \text{NLik}(p_\theta) + \text{KL}(q \parallel p_\theta)$$

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

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

SUMMARY



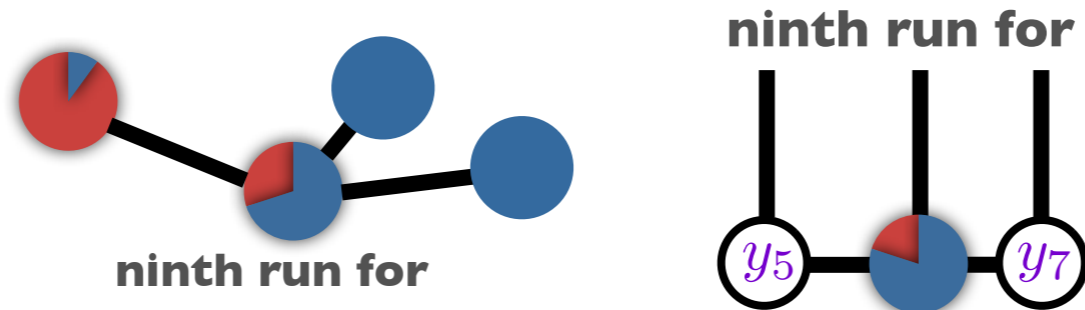
$$\mathcal{J}(q, p_\theta) = \text{Lap}(q) + \text{NLik}(p_\theta) + \text{KL}(q \parallel p_\theta)$$

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

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

Theorem:
Converges to a local
optimum of
 $\mathcal{J}(q, p_\theta)$

SUMMARY



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

any convex, differentiable $g(m)$

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

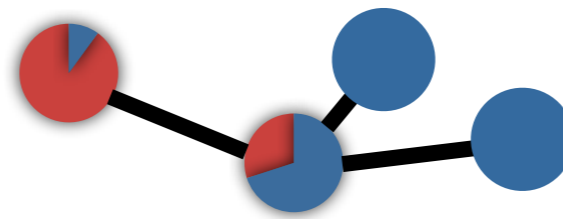
E-step: $q_{\mathbf{y}}^{i'} = \frac{1}{Z_q(\mathbf{x}^i)} q_{\mathbf{y}}^i \exp \left[-\eta \frac{\partial \mathcal{J}(q, p_\theta)}{\partial q_{\mathbf{y}}^i} \right]$

Theorem:

Converges to a local optimum of

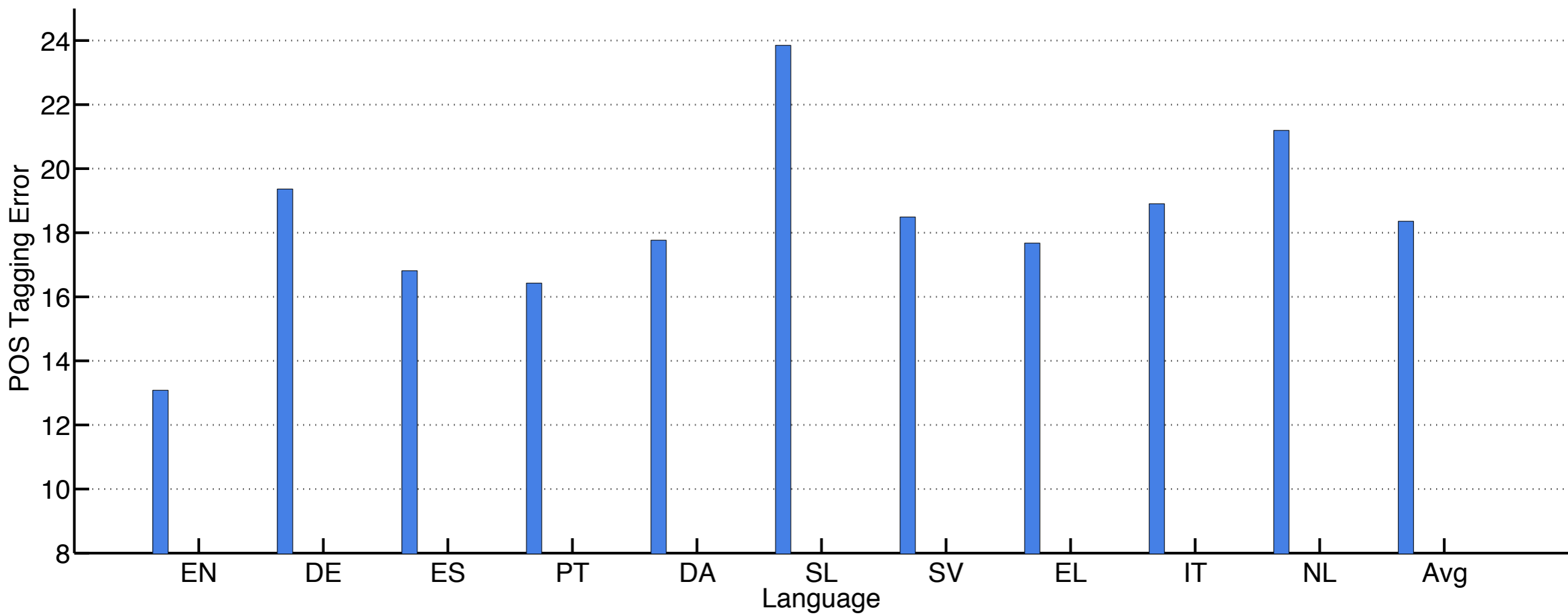
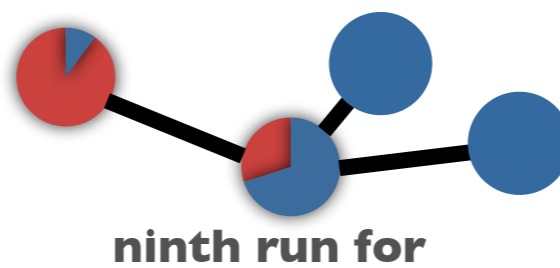
$$\mathcal{J}(q, p_\theta)$$

■ graph-propagation



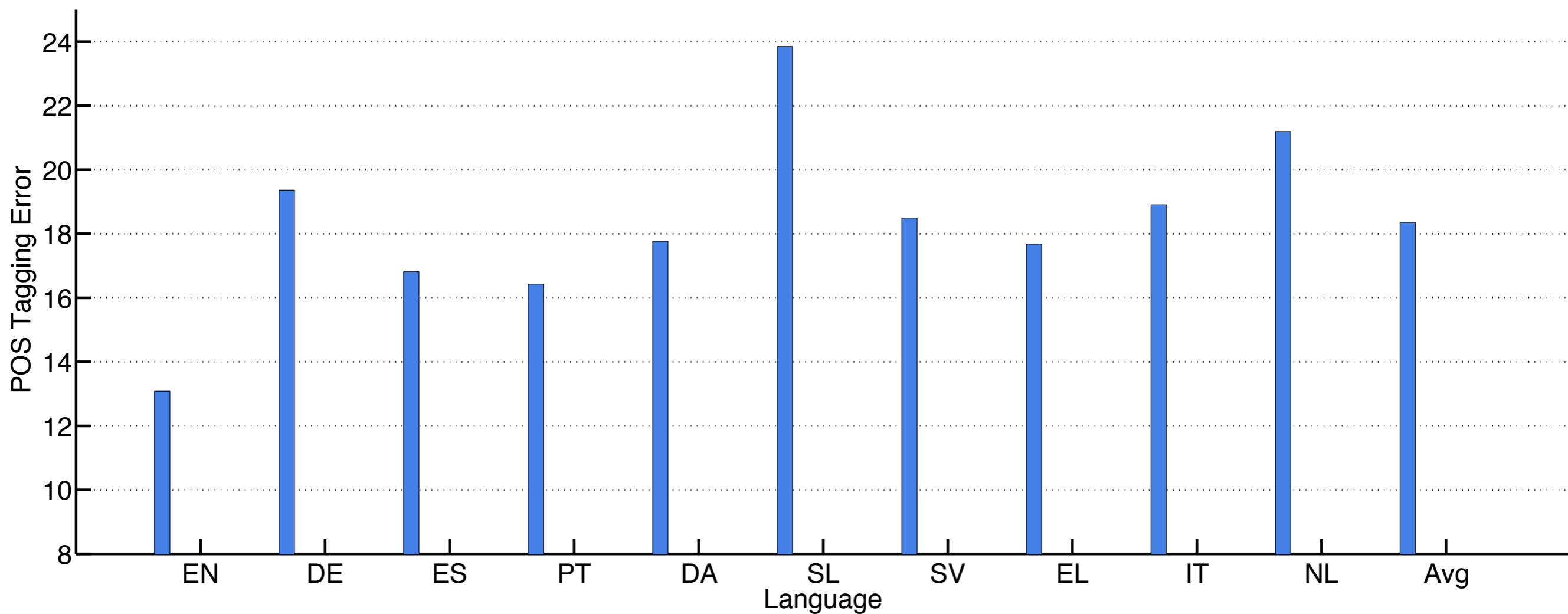
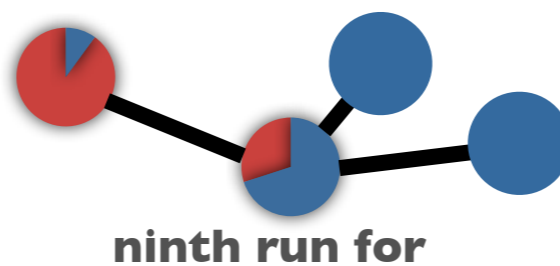
ninth run for

graph-propagation



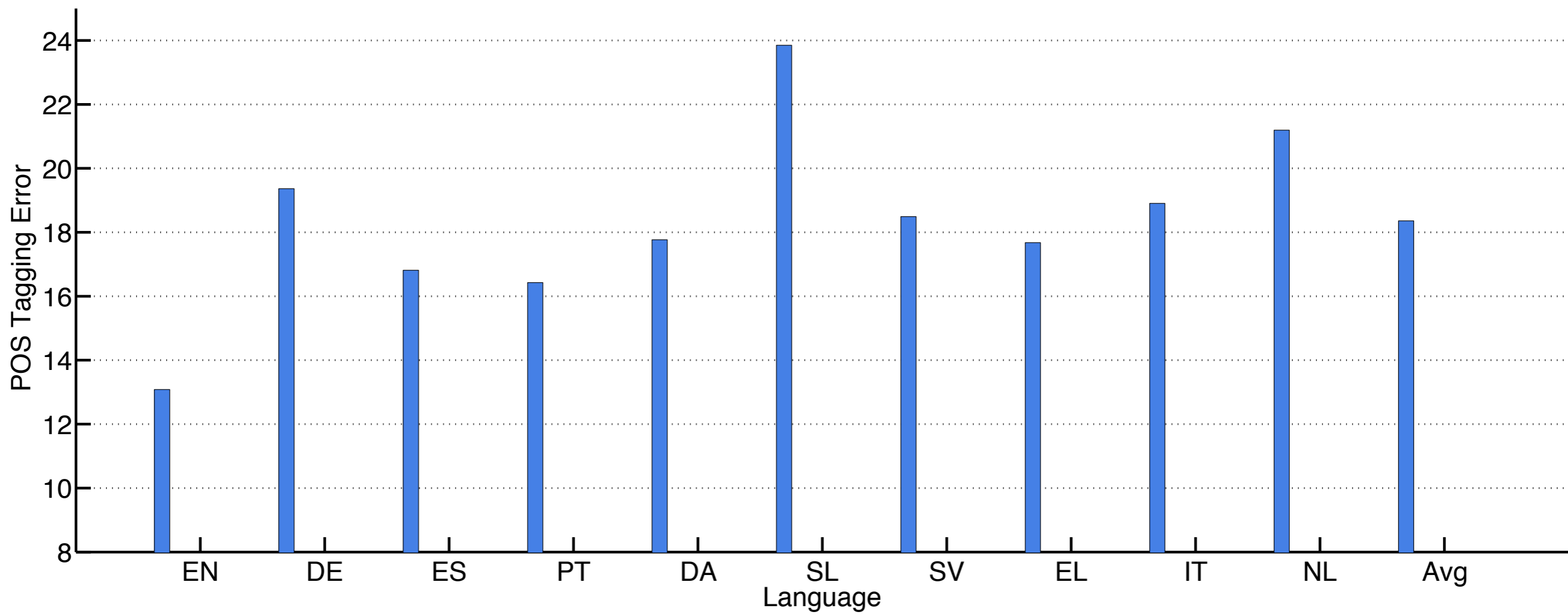
GP

graph-propagation



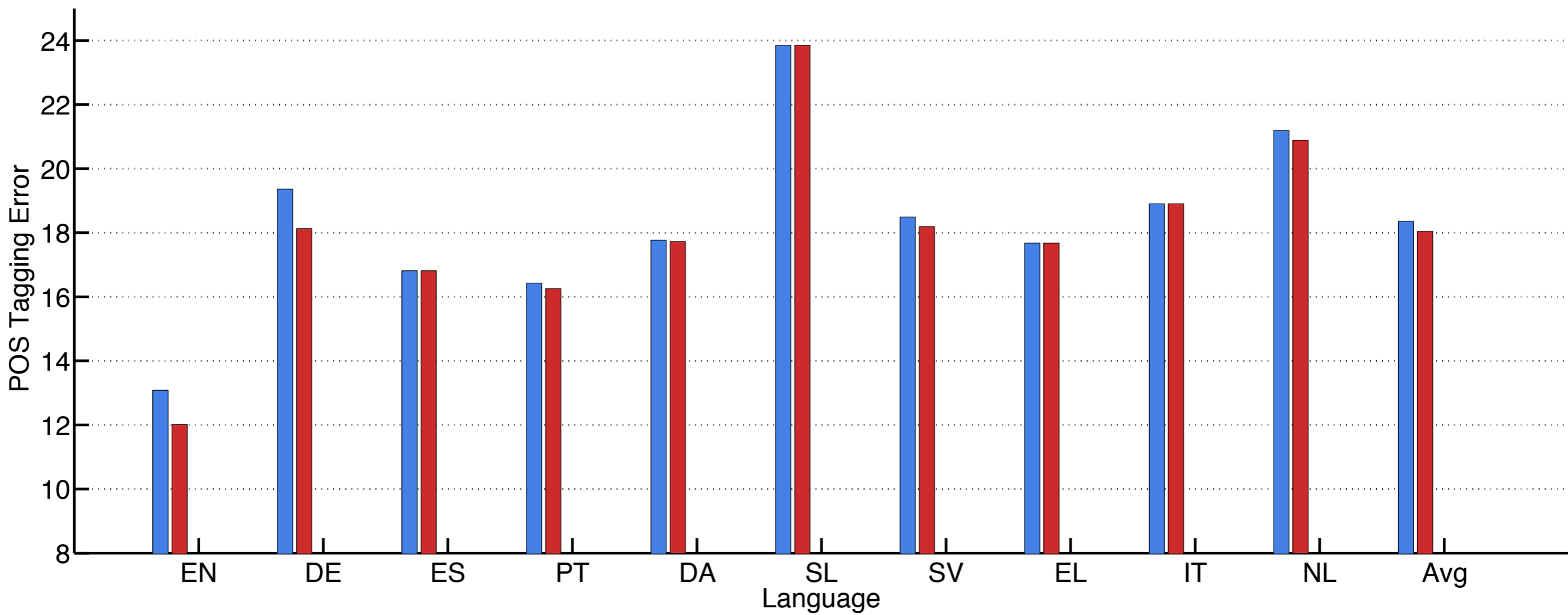
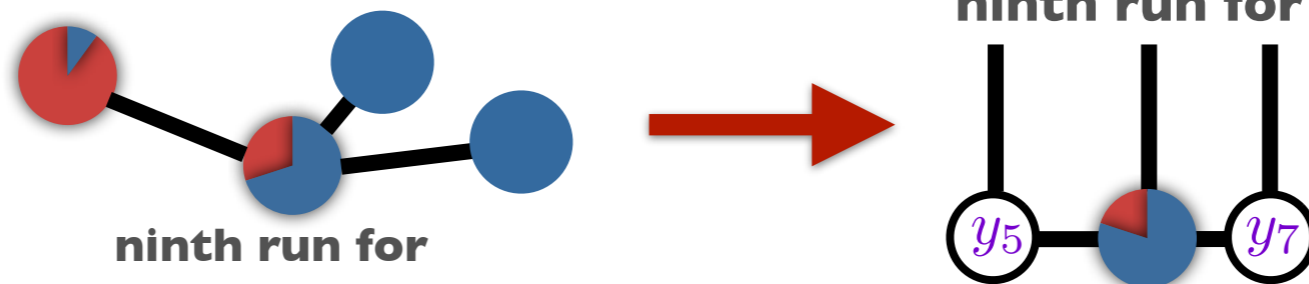
GP

100 labeled sentences



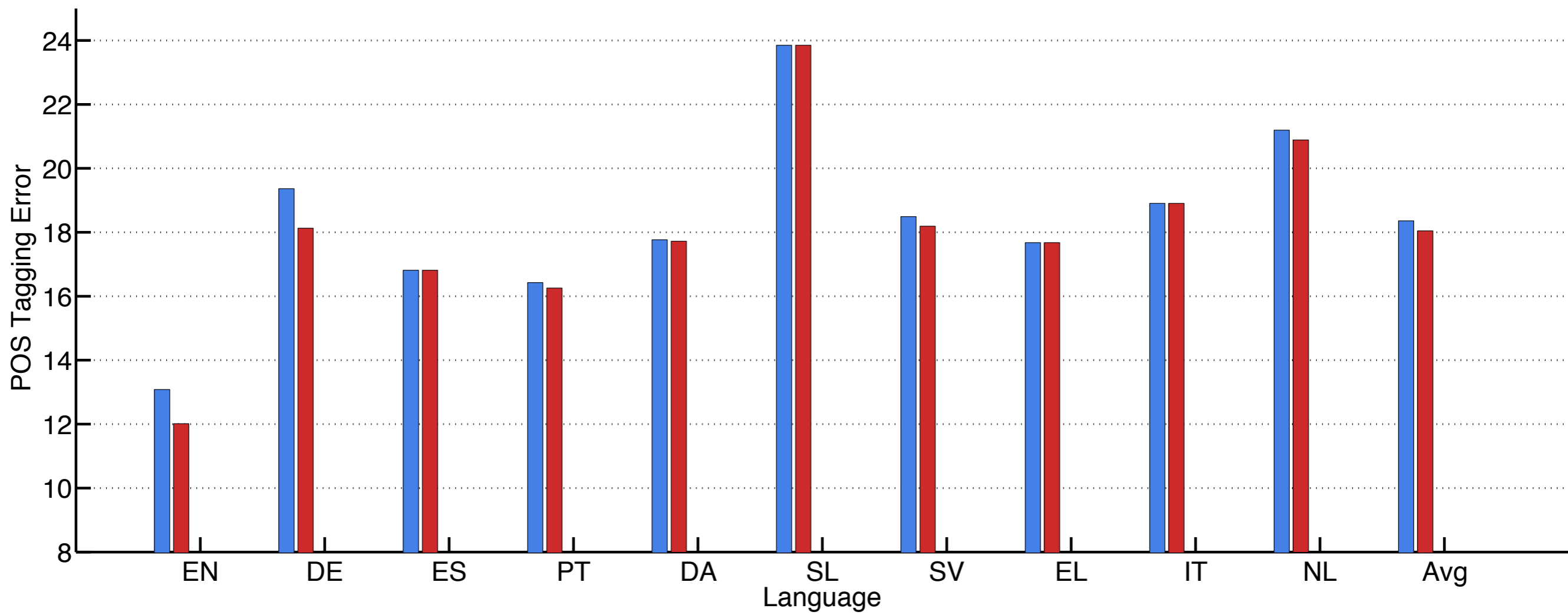
 **GP**

GP → CRF




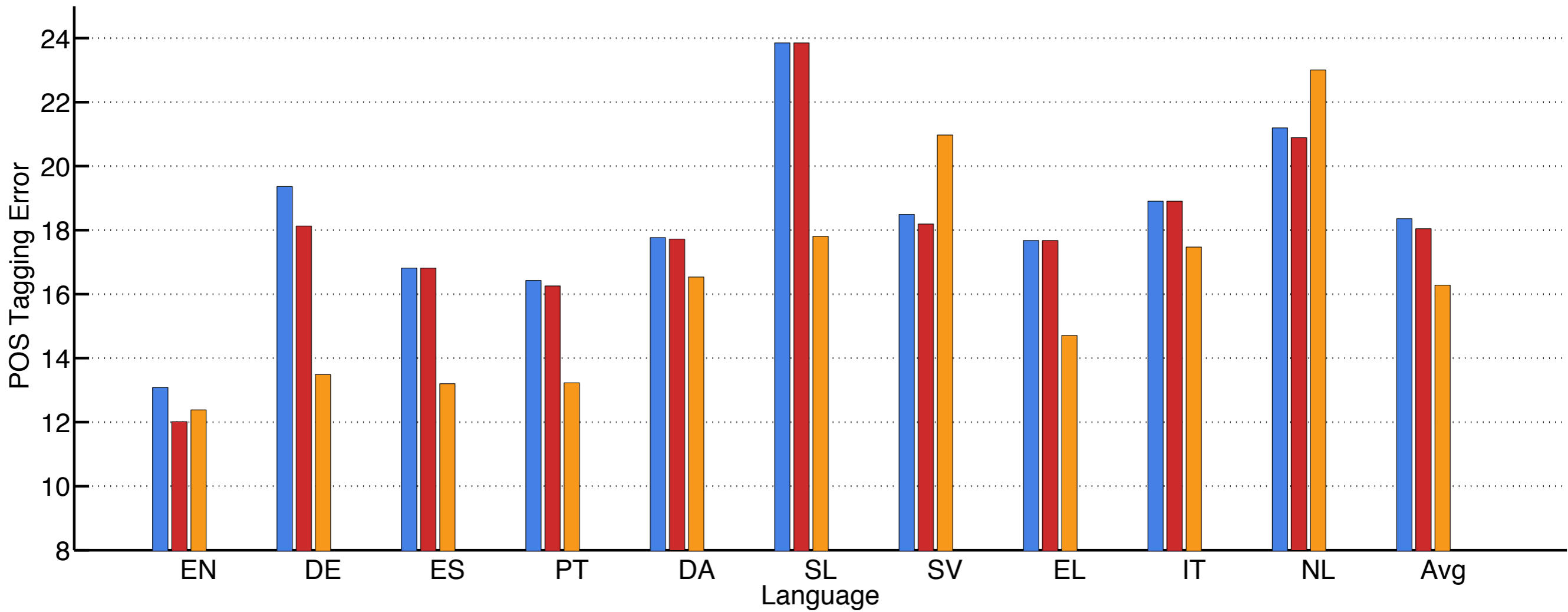
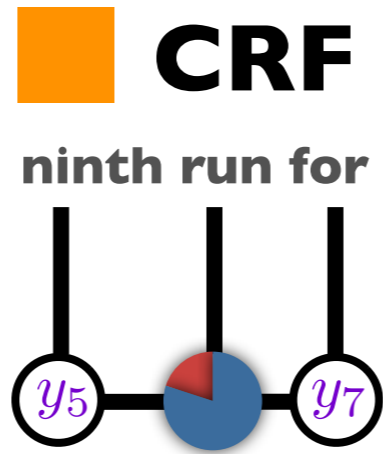
GP

GP → CRF

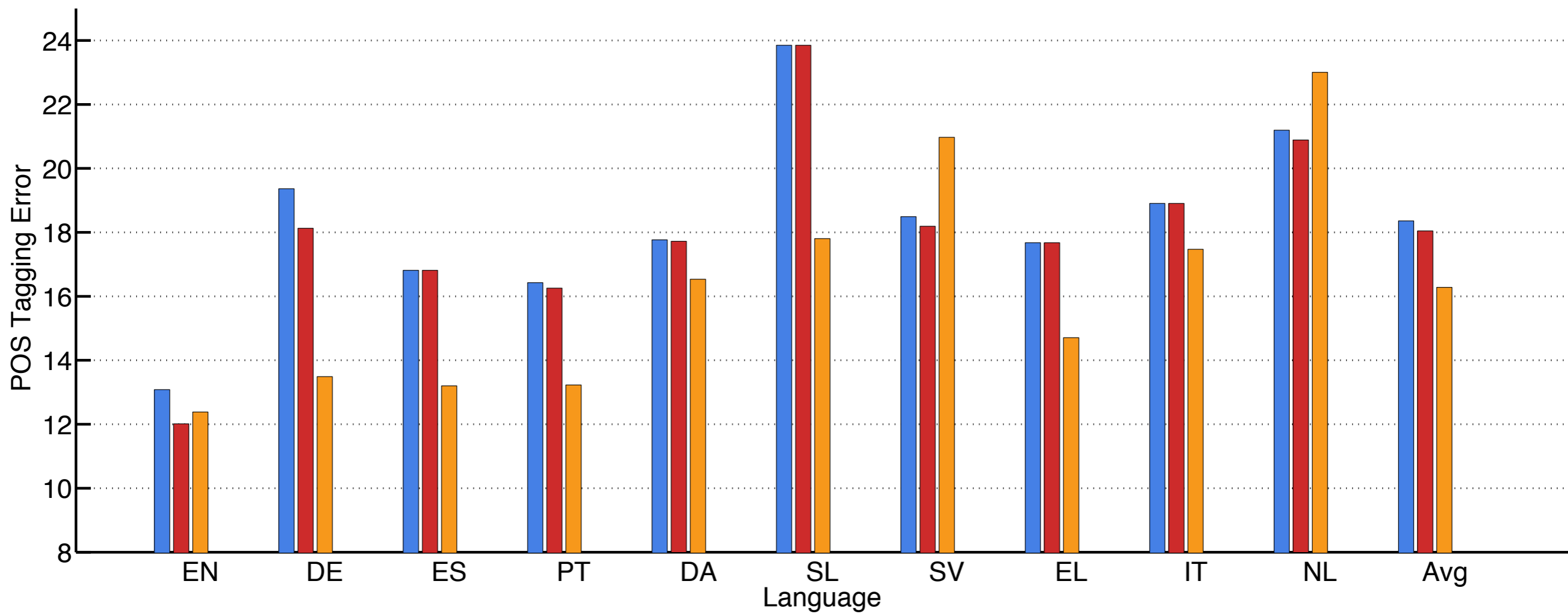


 **GP**

 **GP → CRF**



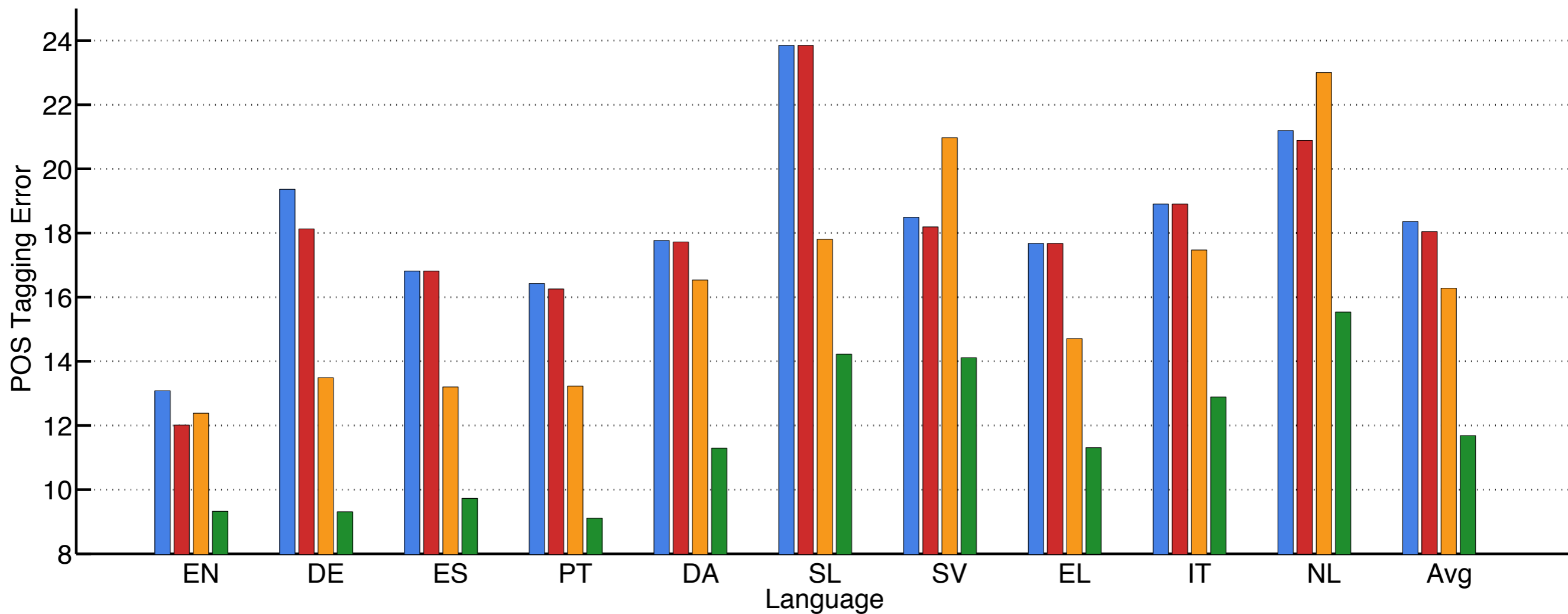
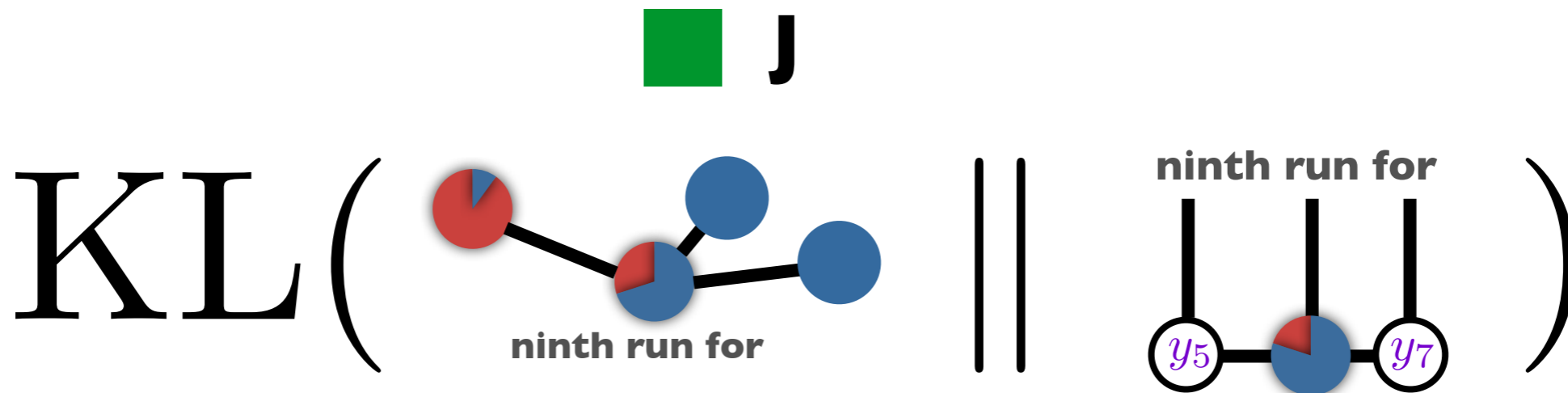
GP **GP → CRF** **CRF**



GP

GP → CRF

CRF

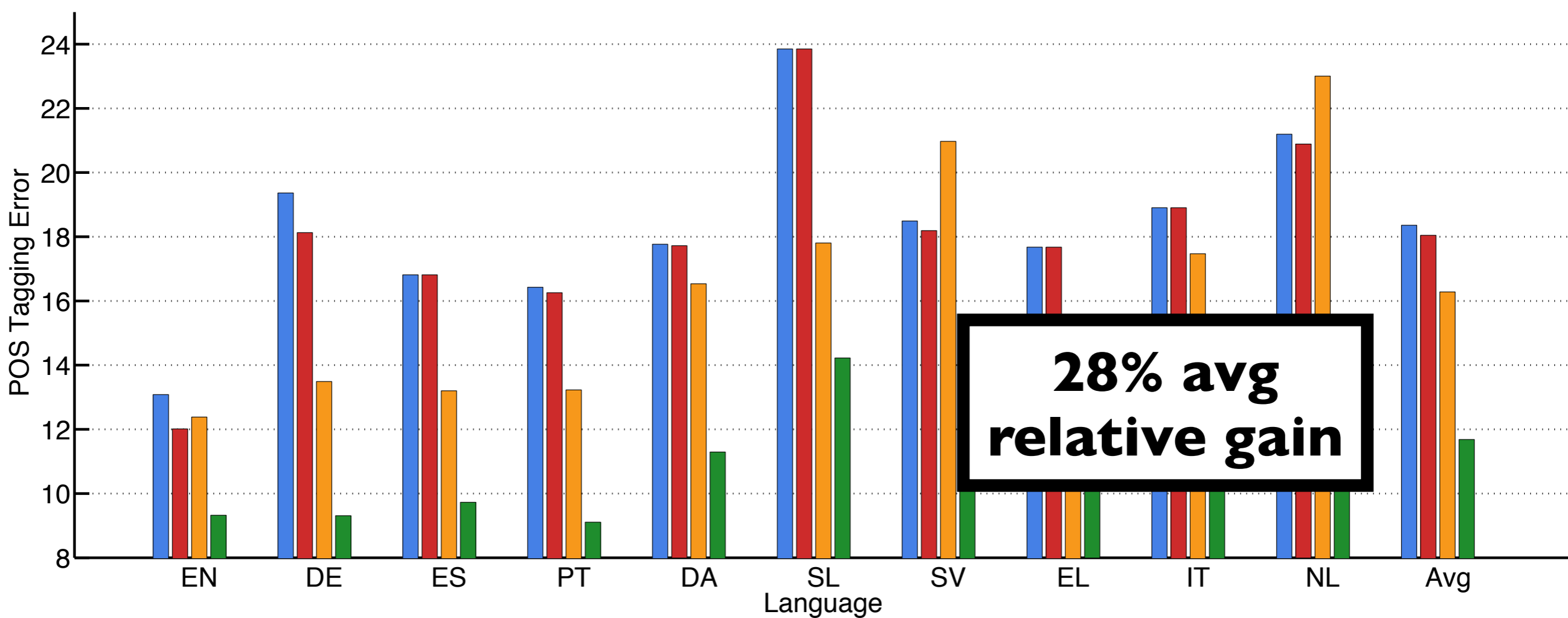
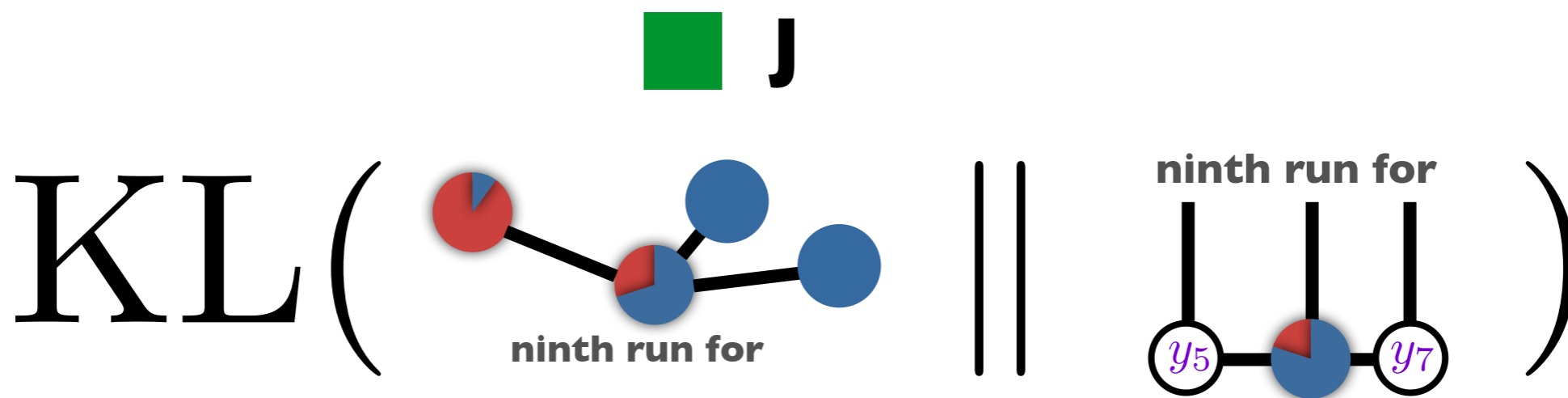


GP

GP → CRF

CRF

J



■ **GP**
■ **GP → CRF**
■ **CRF**
■ **J**

QUESTIONS?

QUESTIONS?

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