

Semi-Supervised Learning of Sequence Models via Method of Moments

EMNLP - Empirical Methods for Natural Language Processing
November 1-6, 2016 Austin, Texas



Zita Marinho

IST, University of Lisbon
Robotics Institute, CMU

zmarinho@cmu.edu

André F. T. Martins

IT, IST, University of Lisbon
Unbabel

andre.martins@unbabel.com

Shay B. Cohen

School of Informatics
University of Edinburgh

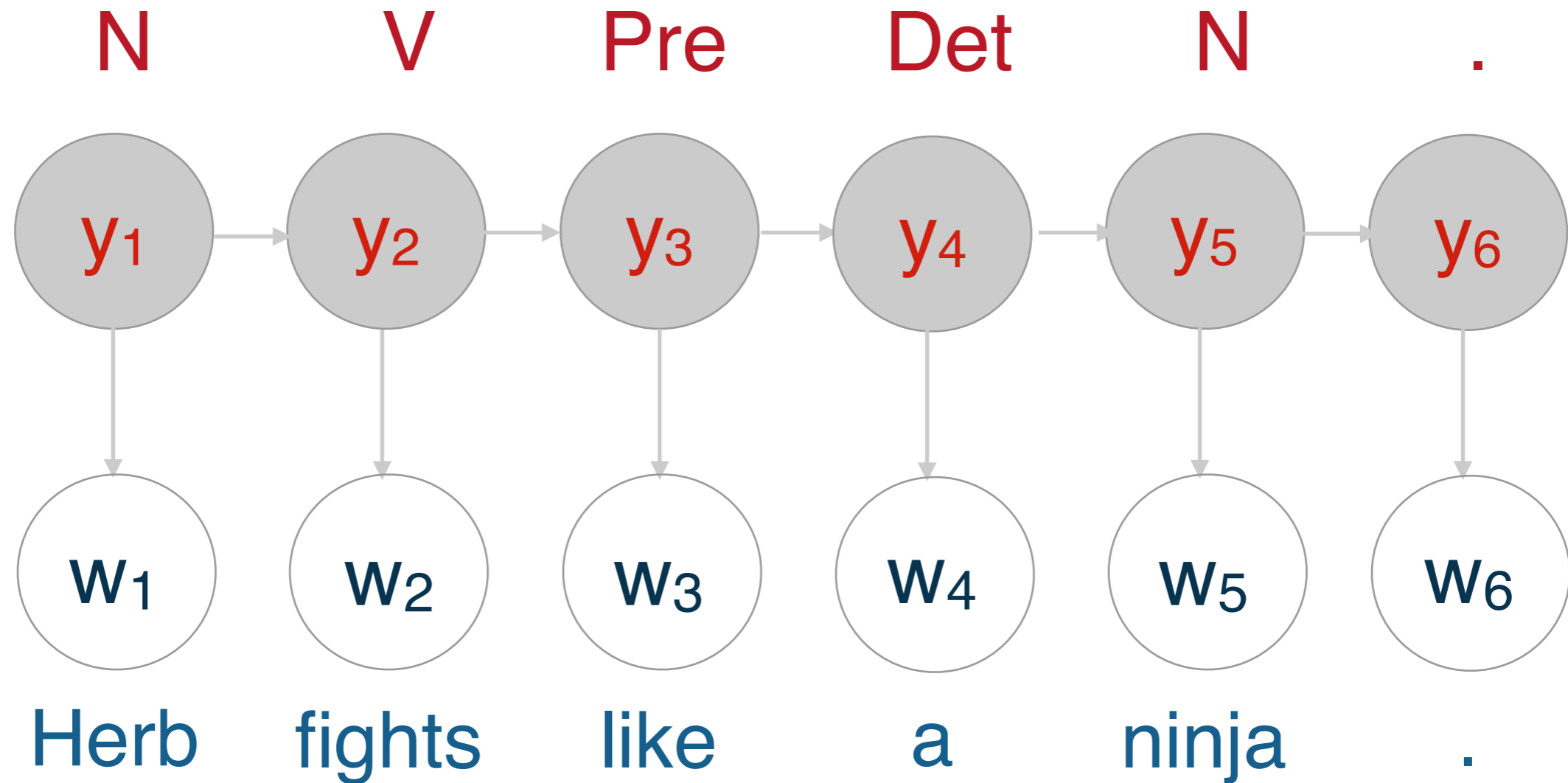
scohen@inf.ed.ac.uk

Noah A. Smith

Computer Science & Eng.
University of Washington

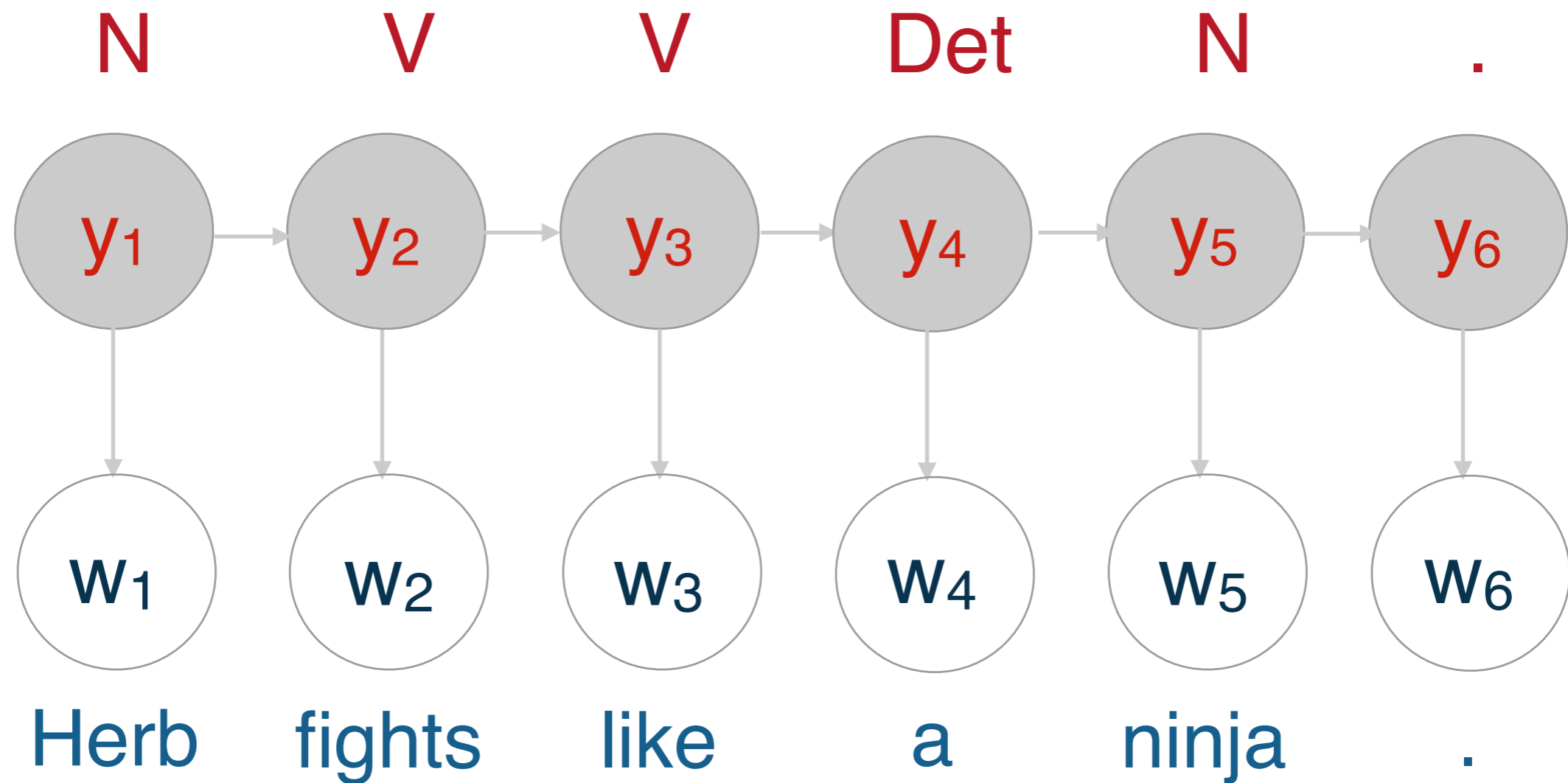
nasmith@cs.washington.edu

Sequence Labeling



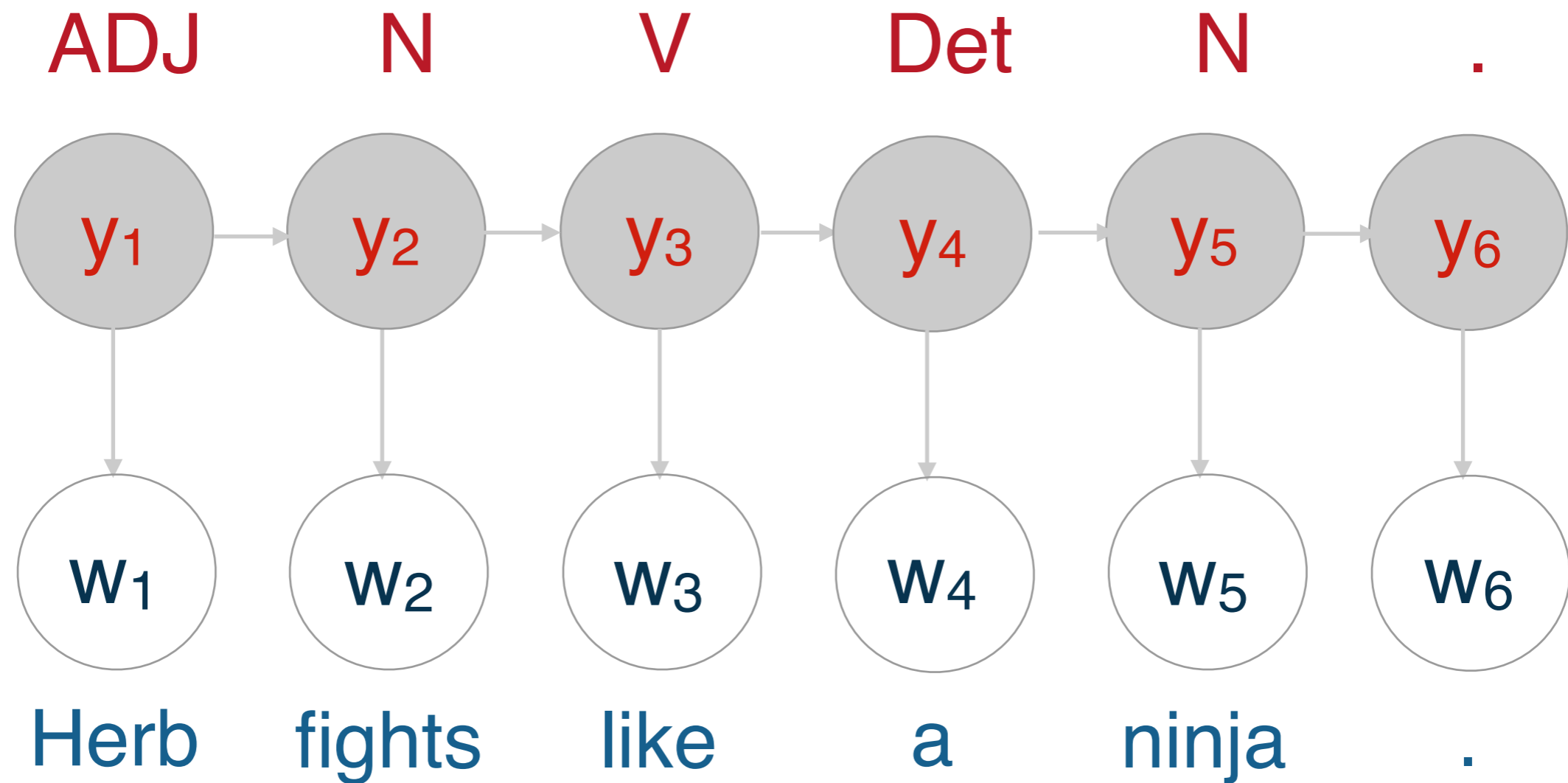
observed data $\{w_1, w_2, w_3, \dots, w_6\}$
labels $\{y_1, y_2, y_3, \dots, y_6\}$

Sequence Labeling



observed data $\{w_1, w_2, w_3, \dots, w_6\}$
labels $\{y_1, y_2, y_3, \dots, y_6\}$

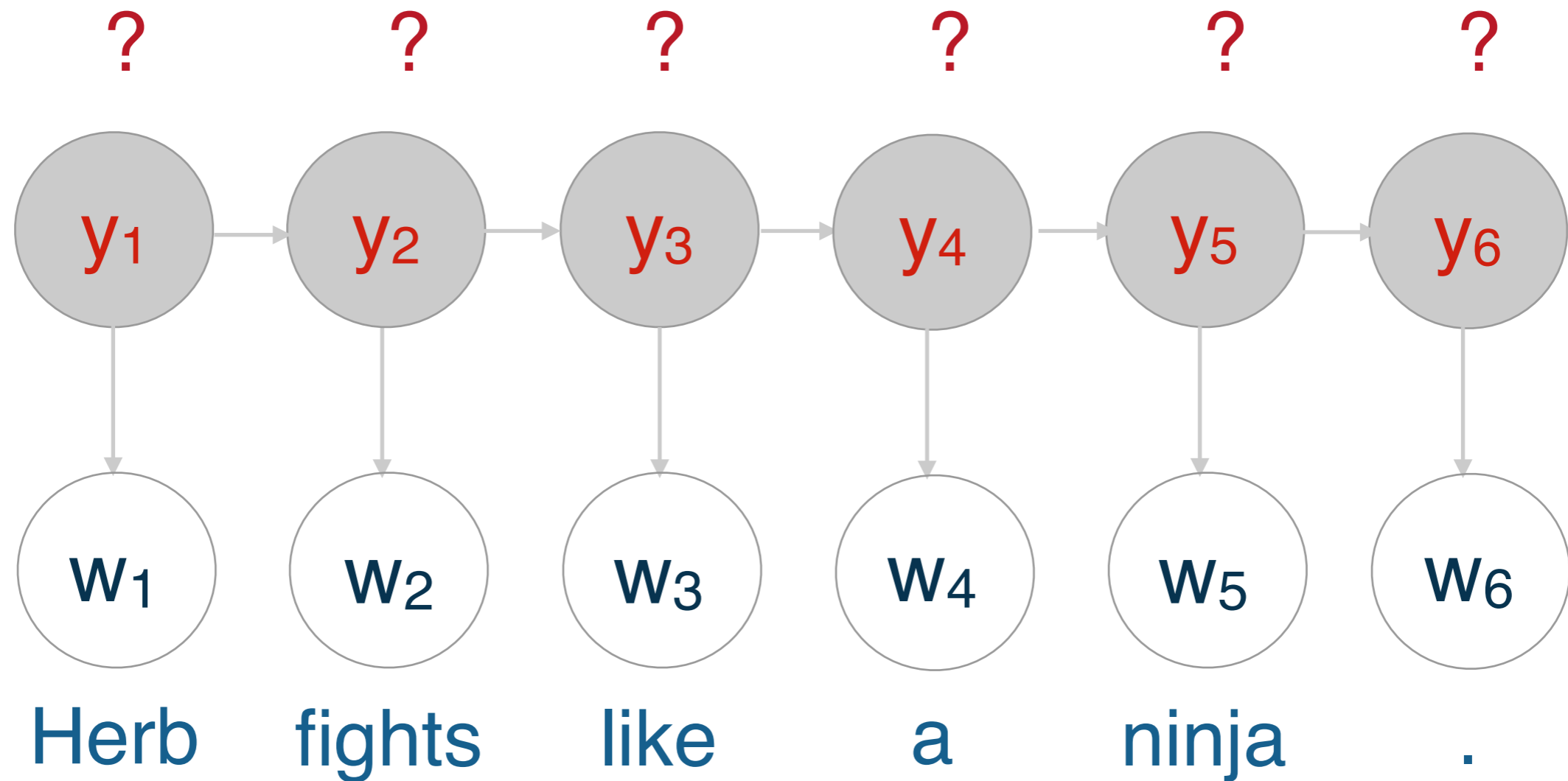
Sequence Labeling



observed data $\{w_1, w_2, w_3, \dots, w_6\}$
labels $\{y_1, y_2, y_3, \dots, y_6\}$

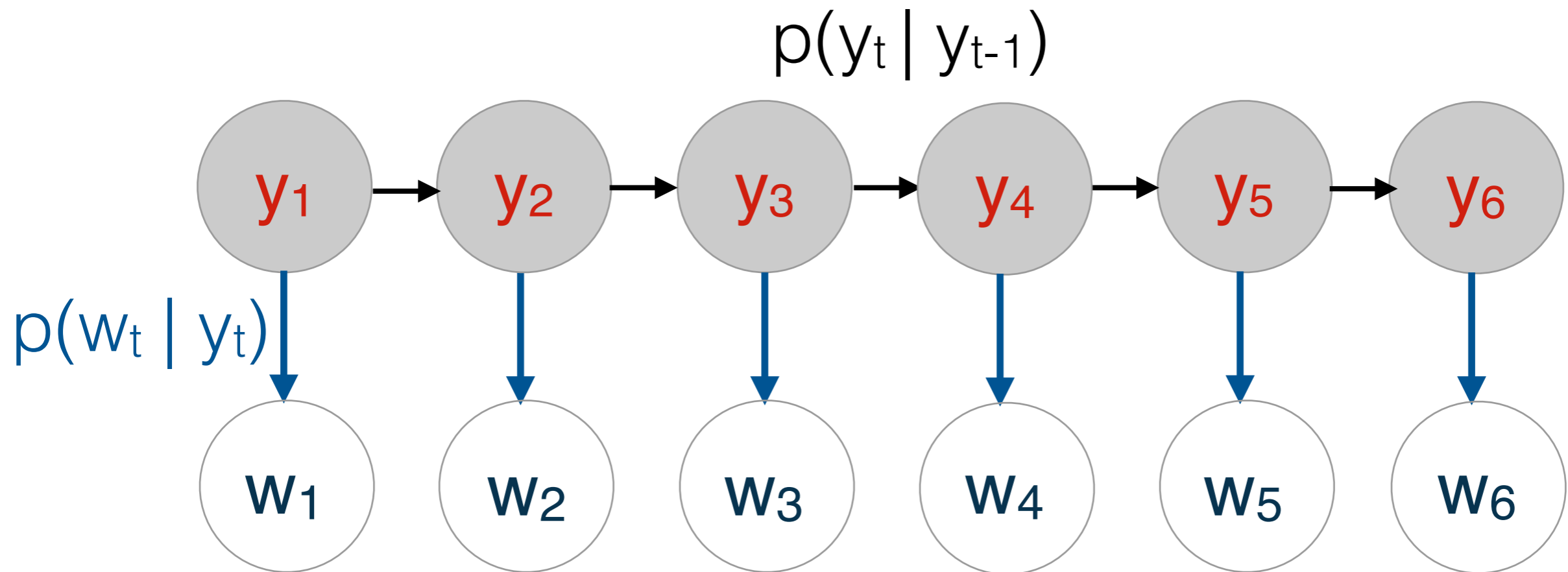
Sequence Labeling

K^6 possible assignments



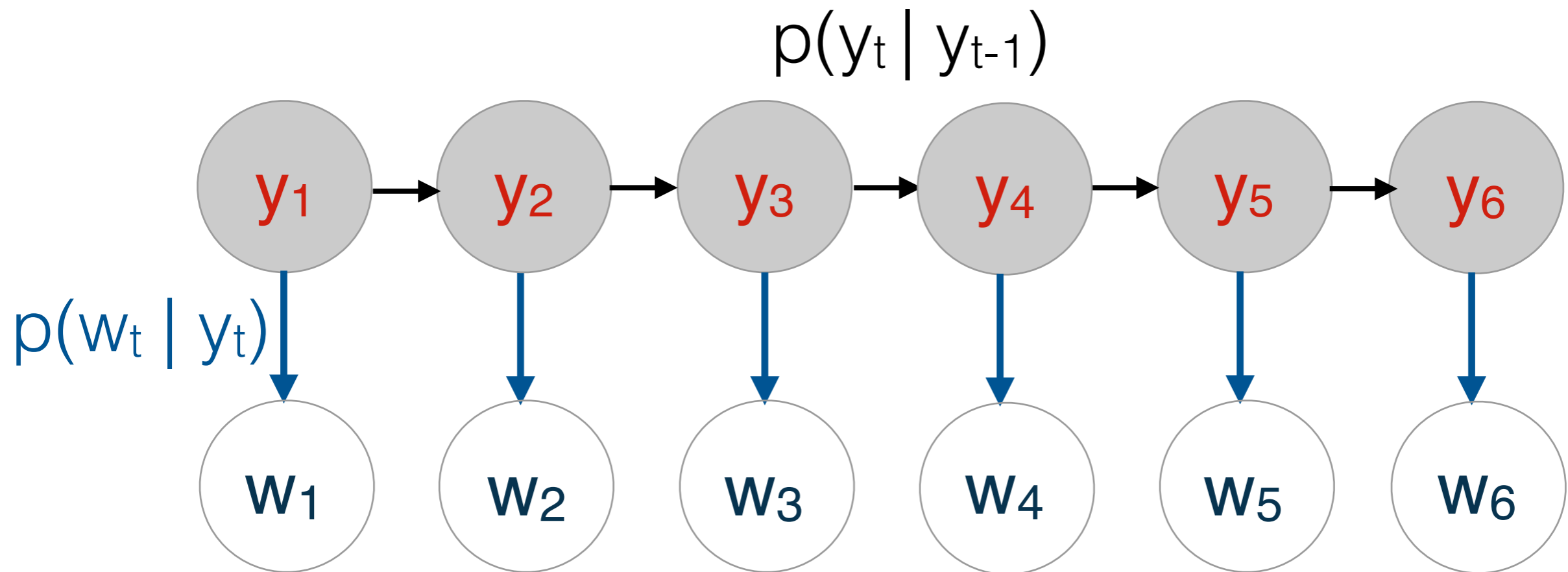
observed data $\{w_1, w_2, w_3, \dots, w_6\}$
labels $\{y_1, y_2, y_3, \dots, y_6\}$

Learn parameters?



- supervised learning
- unsupervised/semi-supervised (this talk)

Learn parameters?



- supervised learning
- unsupervised/semi-supervised (this talk)
- model can be extended to include features

Berg-Kirkpatrick, et al, Painless unsupervised learning with features. NAACL HLT, 2010.

Maximum Likelihood estimation (MLE)

Method of Moments estimation (MoM)

- exact inference is hard → **computationally efficient**
- EM sensitive to local optima (depends on initialization) → **no local optima**
- EM expensive in large datasets (several inference passes) → **one pass over data**

Hidden Markov Model

via Maximum Likelihood Estimation

via Method of Moments

	MLE HMM	MLE feature HMM	MoM HMM	MoM feature HMM
semi-supervised learning	✓	✓	?	?
unsupervised learning	✓	✓	✓	?

Shay B. Cohen, Karl Stratos, Michael Collins, Dean P. Foster and Lyle Ungar, *Spectral Learning of Latent-Variable PCFGs: Algorithms and Sample Complexity*, JMLR 2014

Arora et al., *A Practical Algorithm for Topic Modeling with Provable Guarantees*, ICML 2013

Learning sequence models via MoM

Outline

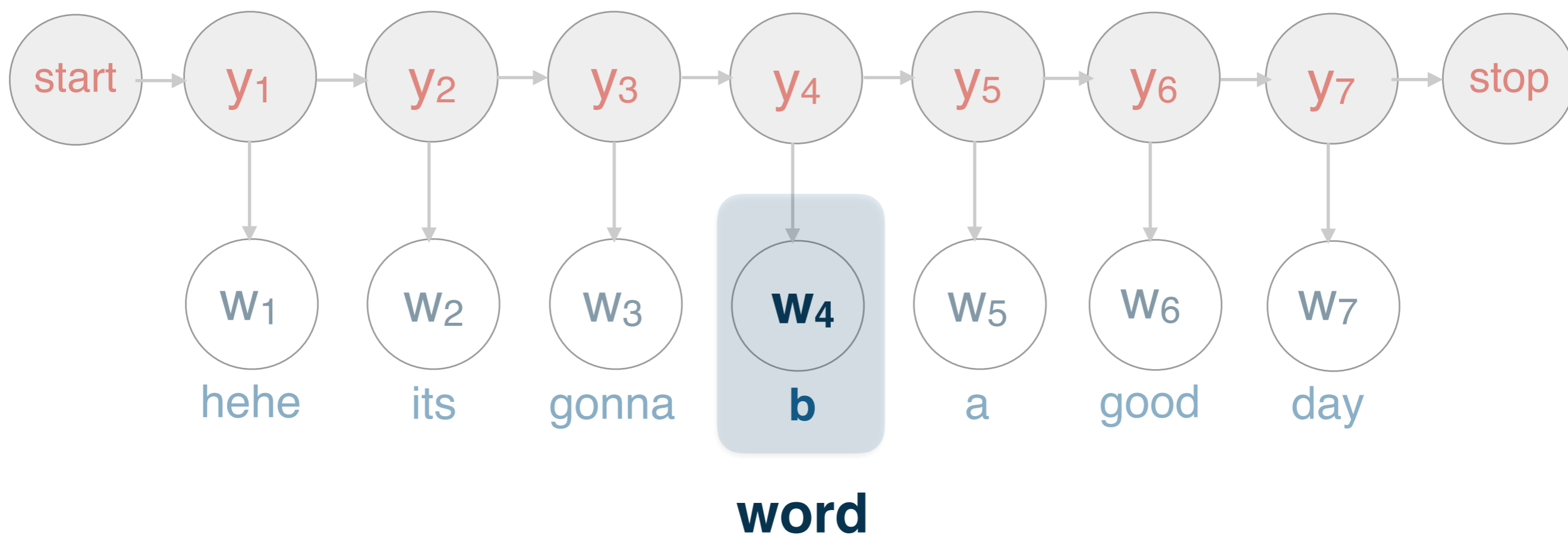
1. Learn HMM models via MoM
2. Solve a QP
3. Extend to feature-based model
4. Experiments



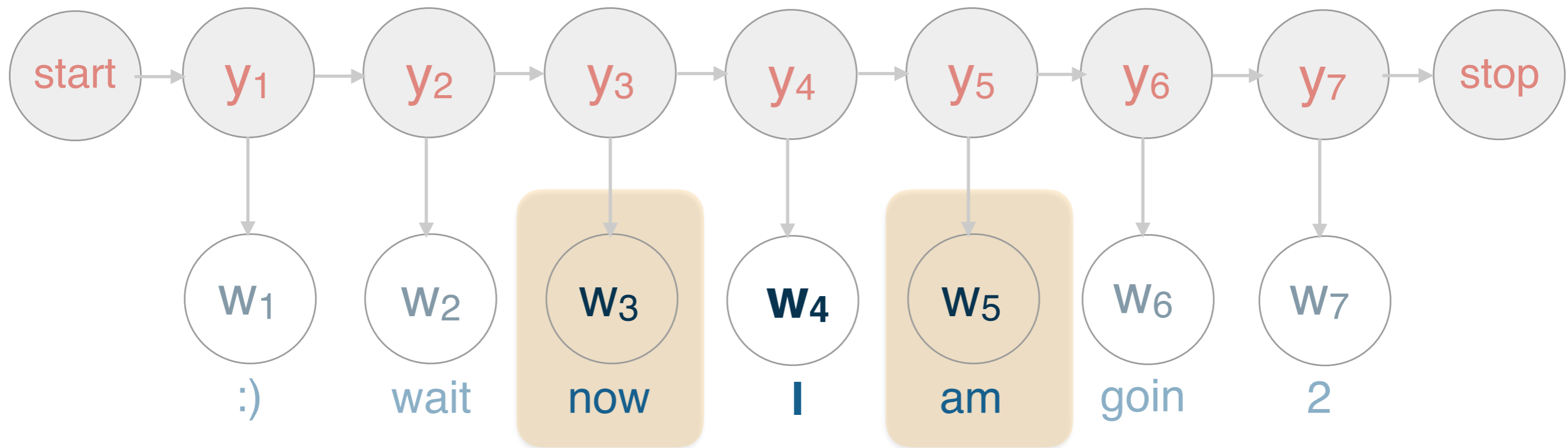
Key insight:

1. **Conditional Independence:**
infer label by looking at context
2. **Anchor Trick:**
learn a proxy for labels with anchors

1. Conditional Independence

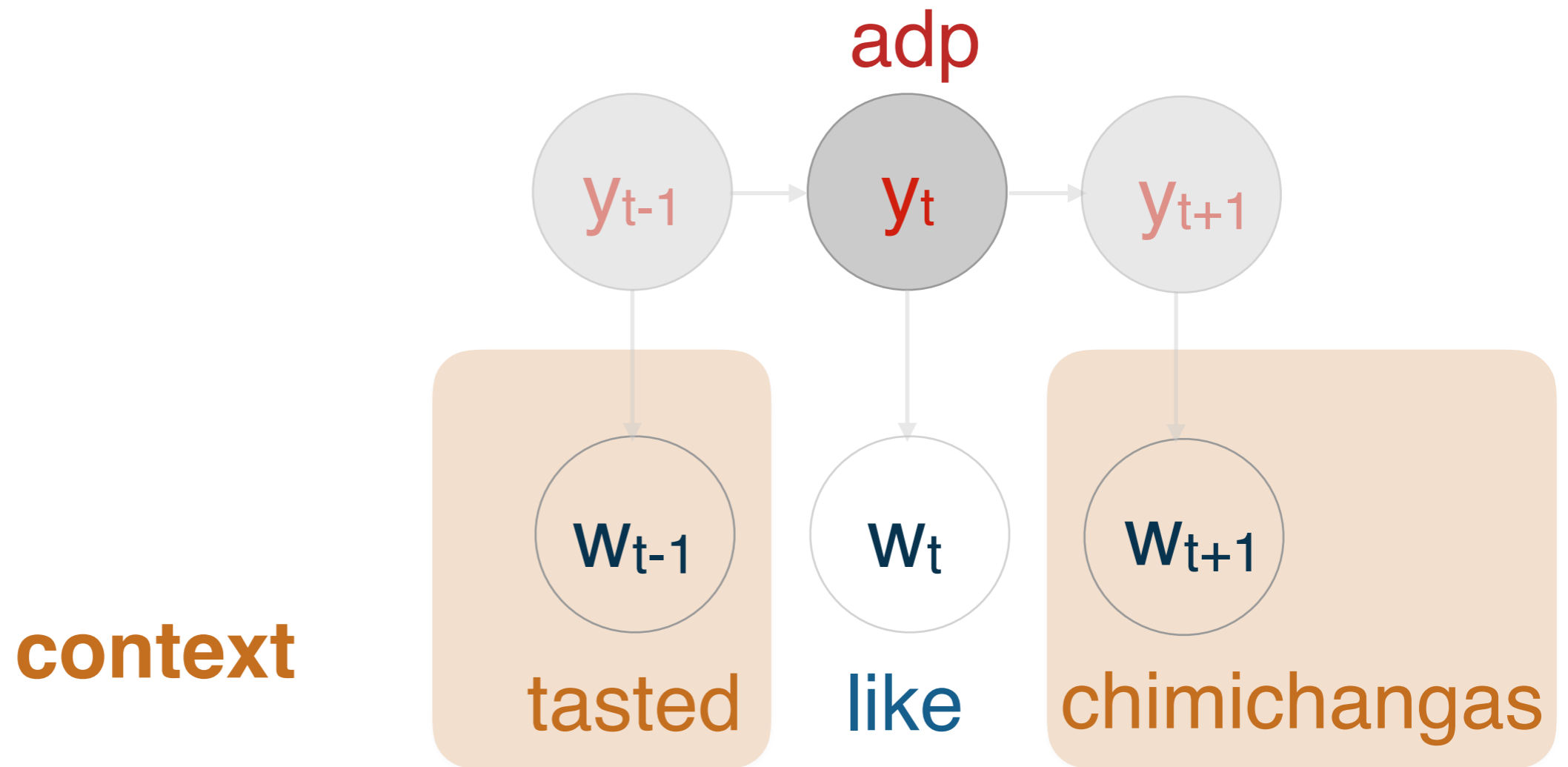


1. Conditional Independence

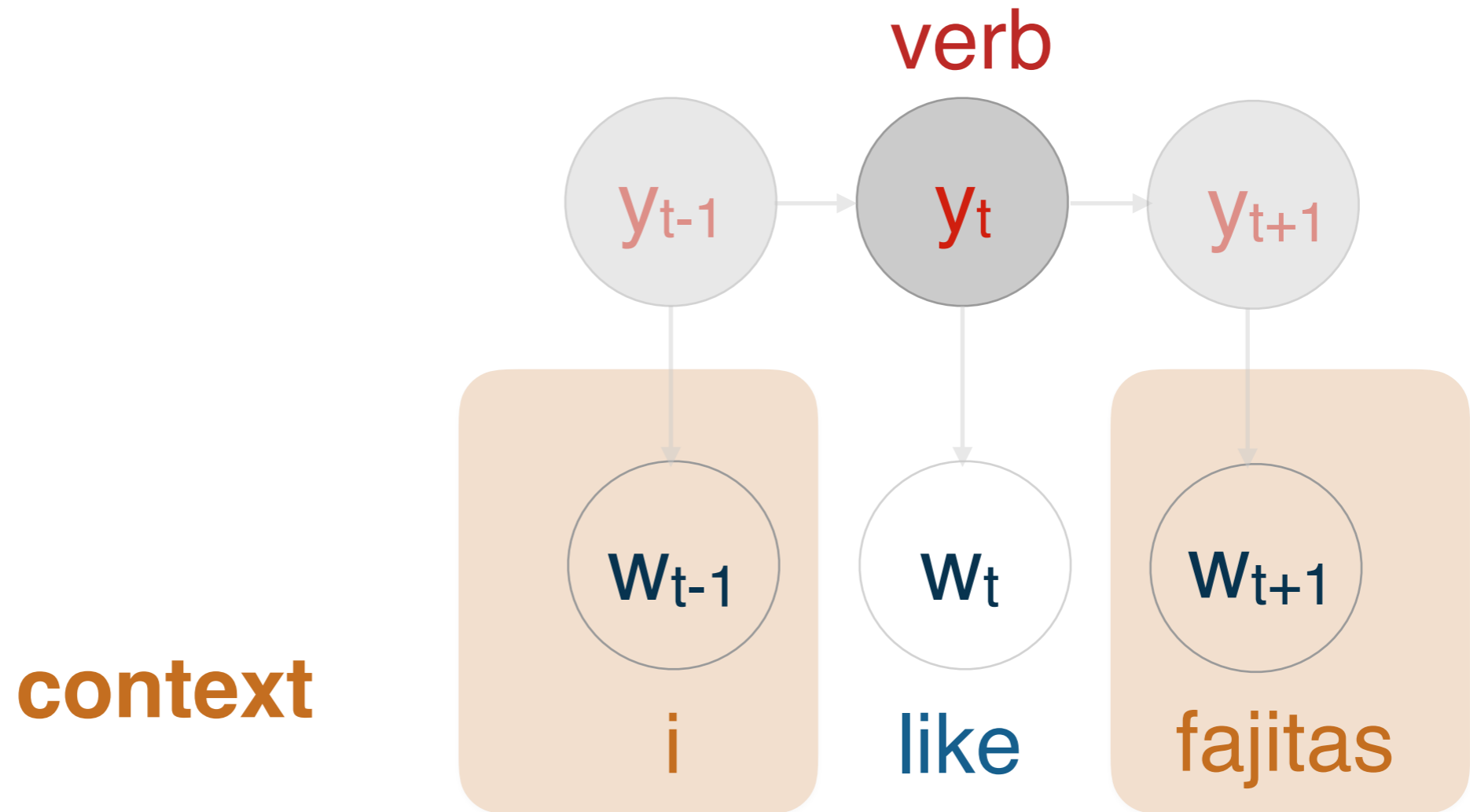


context = { w_{-1} , w_{+1} }

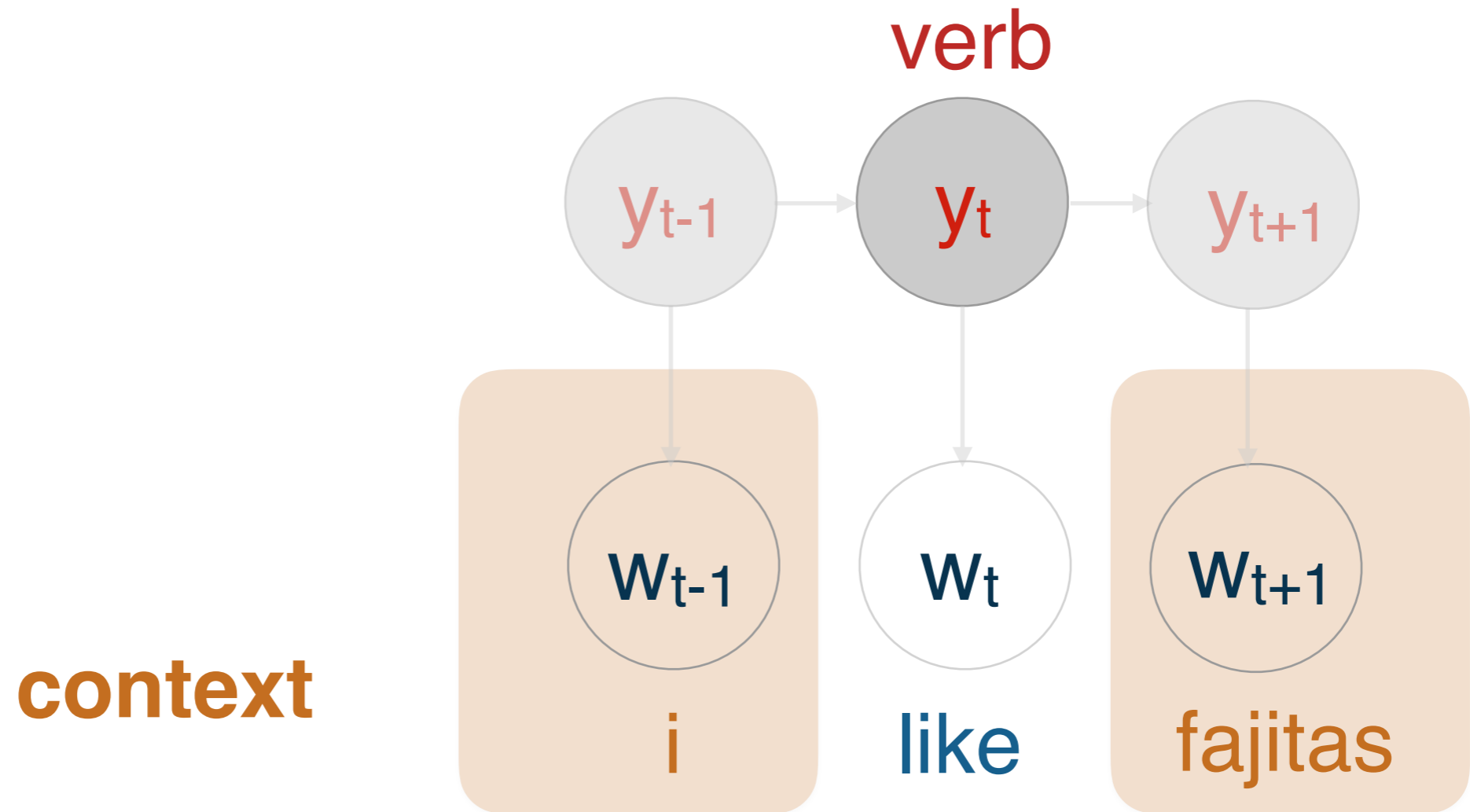
1. Conditional Independence



1. Conditional Independence



1. Conditional Independence

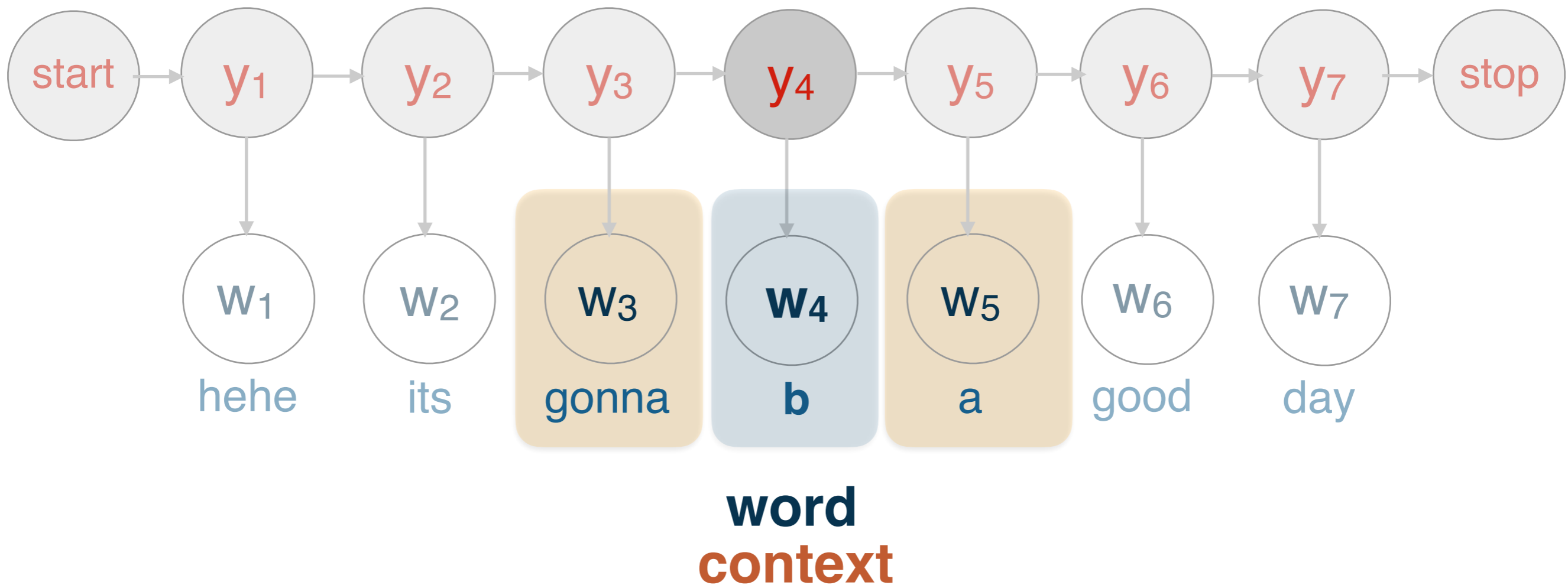


“You shall know a **word** by the **company** it keeps.”

Firth, 1957

1. Conditional Independence

$$\text{word} \perp \text{context} \mid \text{label}$$

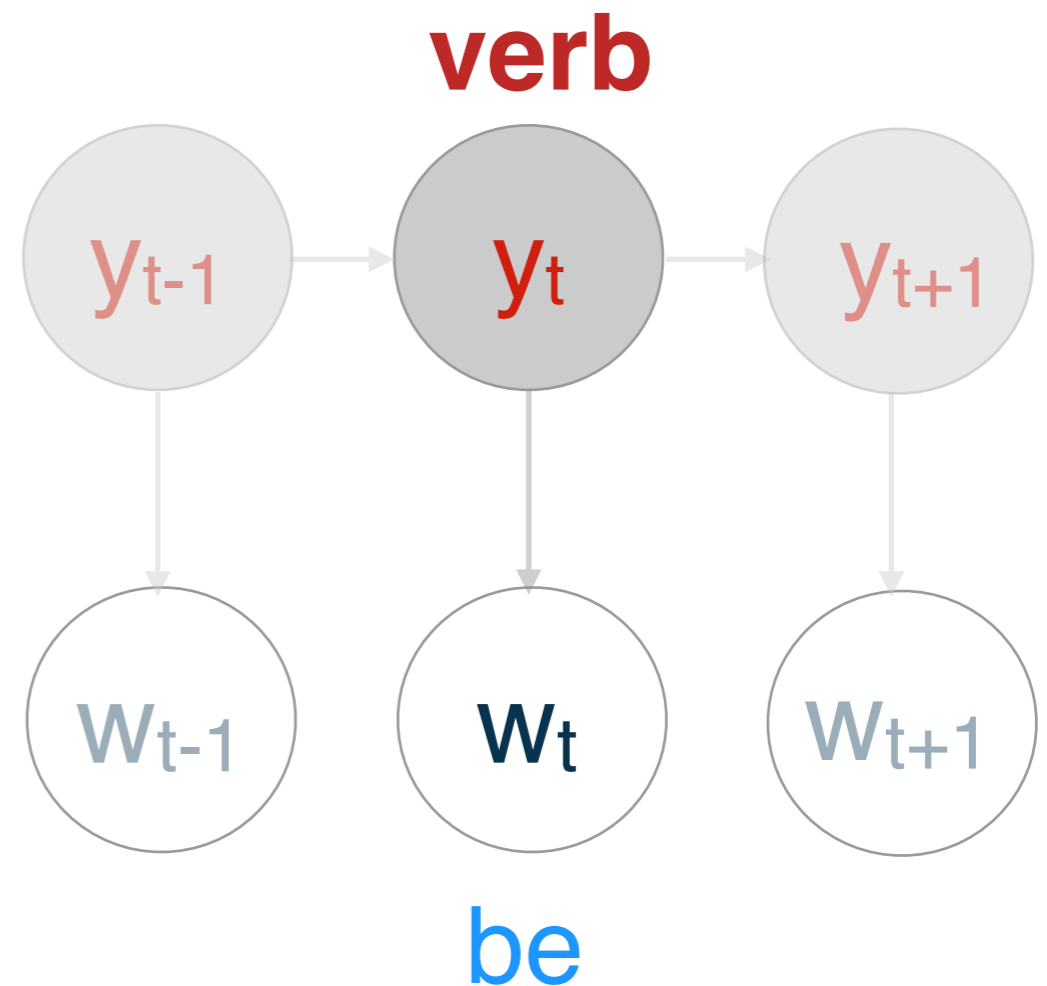


2. Anchor Trick

$$p(\text{verb} \mid \text{be}) = 1$$

$$p(\text{label} \neq \text{verb} \mid \text{be}) = 0$$

all instances of **be** = **verb** ⚓

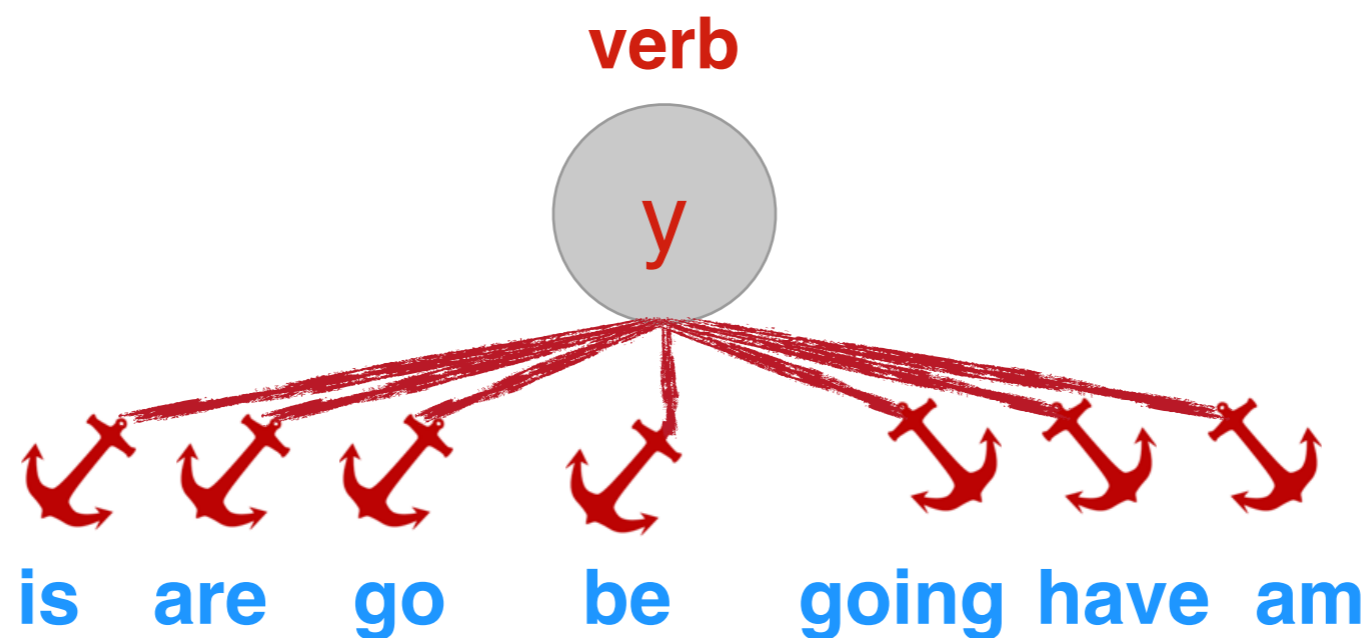


Arora et al., *A Practical Algorithm for Topic Modeling with Provable Guarantees*, ICML 2013

More anchors per label

2. Anchor Trick

verb = **b, be, are, is, am, have, going**



more than 1 anchor word



less biased context estimates

2. Anchor Trick

How to find **anchors**?

- small labeled corpus
- small lexicon

Austin
airport
playground

noun

am,be,is,are
go,
make,made
become

verb

he,it,she

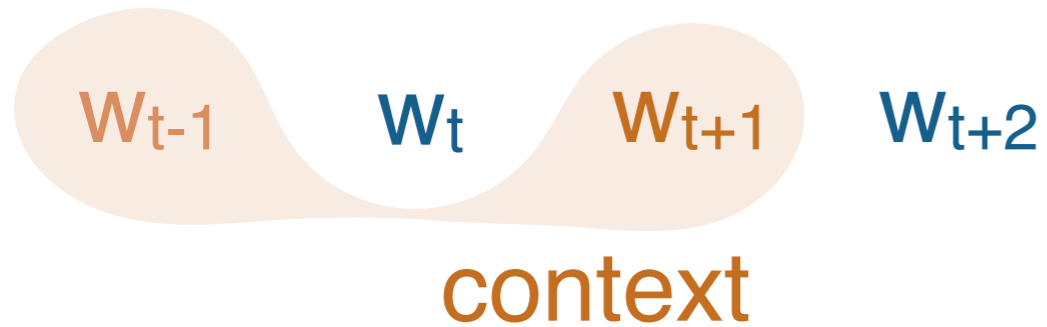
pron

so,on,of

adp

unlabeled

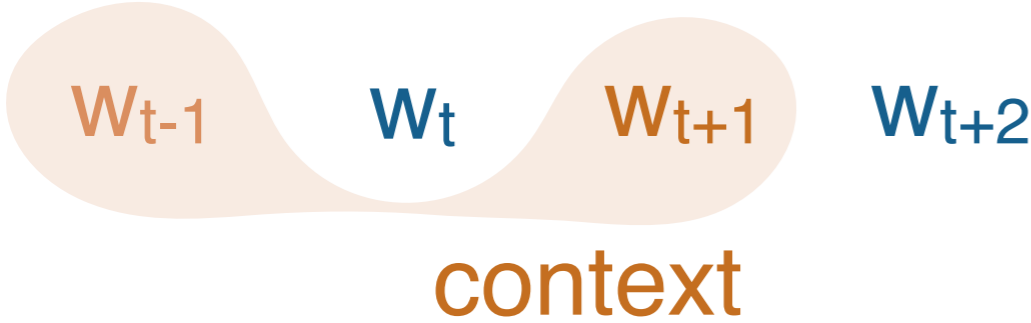
co-occurrences in data



Andrew fights like Jet Li.
Ann sings like me.

eat Fruit like cherry.
Children like ice-cream.

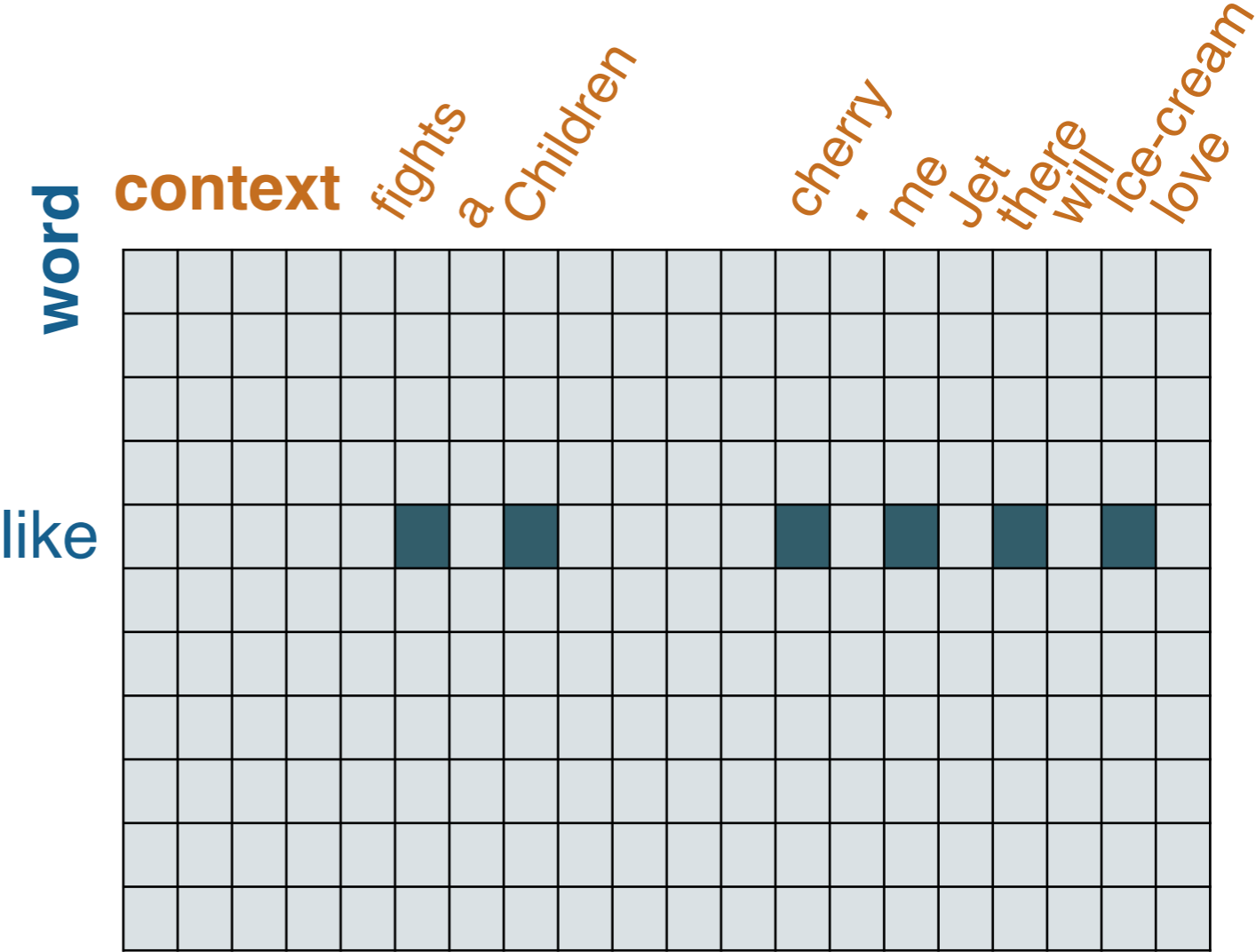
Method of moments



Q $p(\text{context} \mid \text{word})$

Andrew fights like Jet Li.
Ann sings like me.

eat Fruit like cherry.
Children like ice-cream.

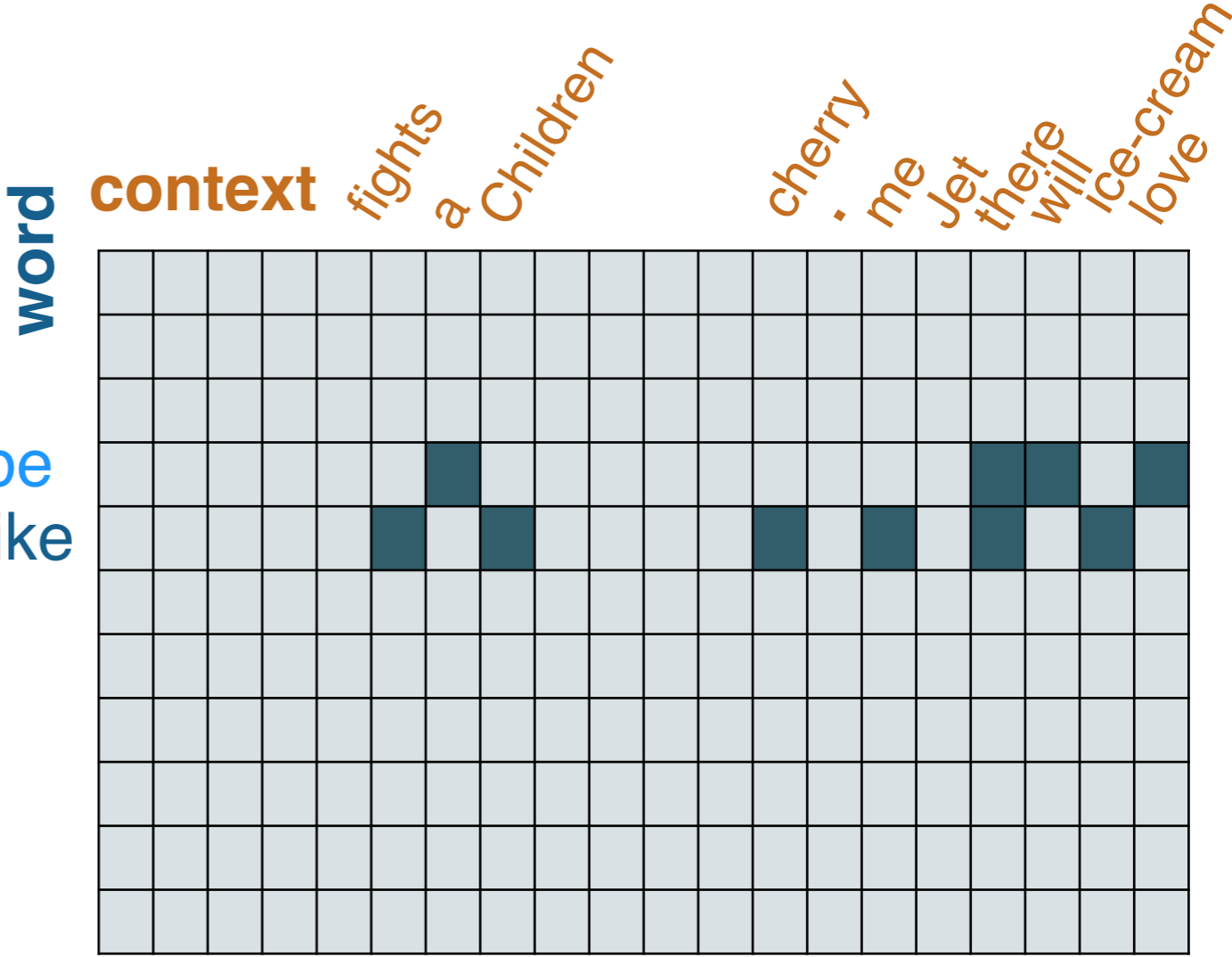


Method of moments

Q

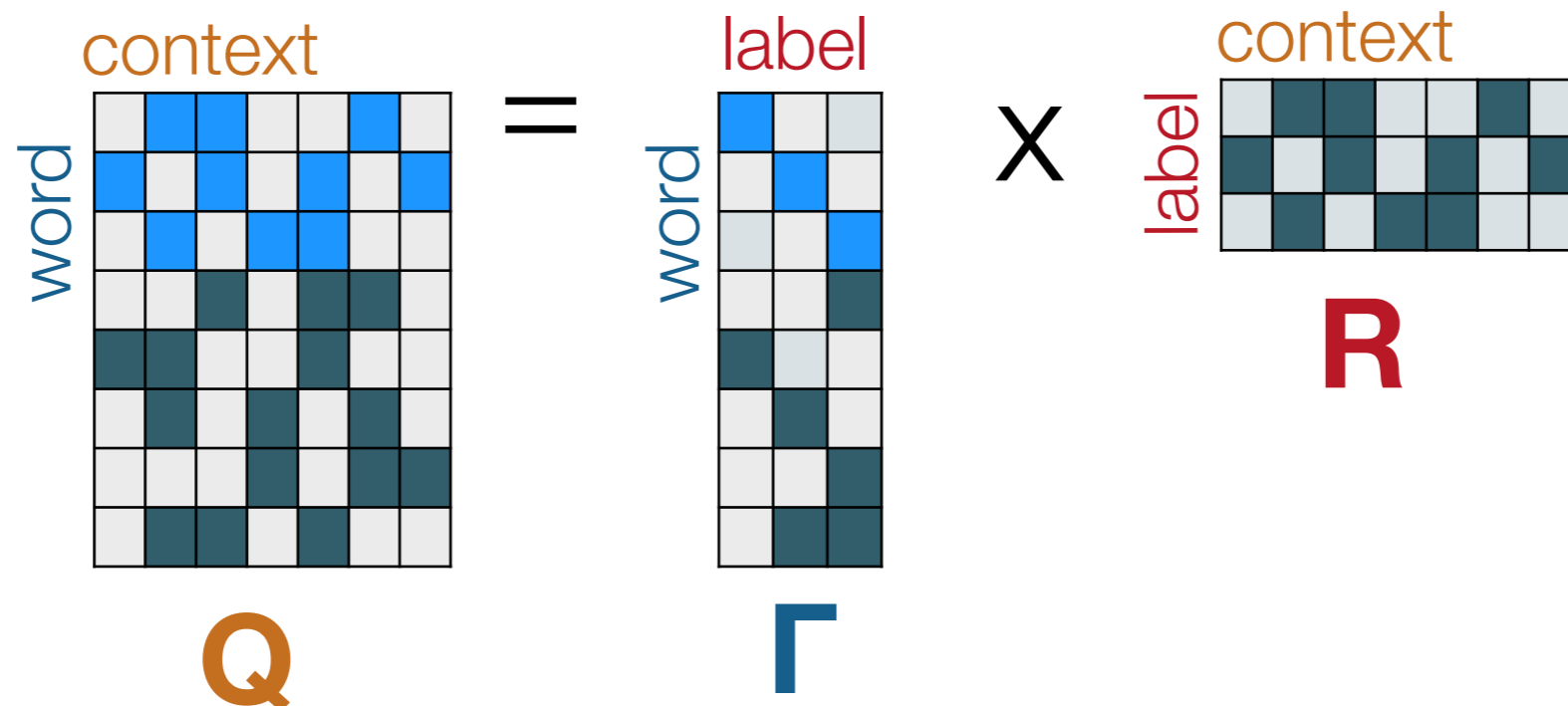
$$p(\text{context} \mid \text{word})$$

Let there be love.
 Bill will be a ninja.



1. Conditional Independence $\text{word} \perp \text{context} \mid \text{label}$

$$p(\text{context} \mid \text{word}) = \sum_{\text{labels}} p(\text{label} \mid \text{word}) p(\text{context} \mid \text{label})$$

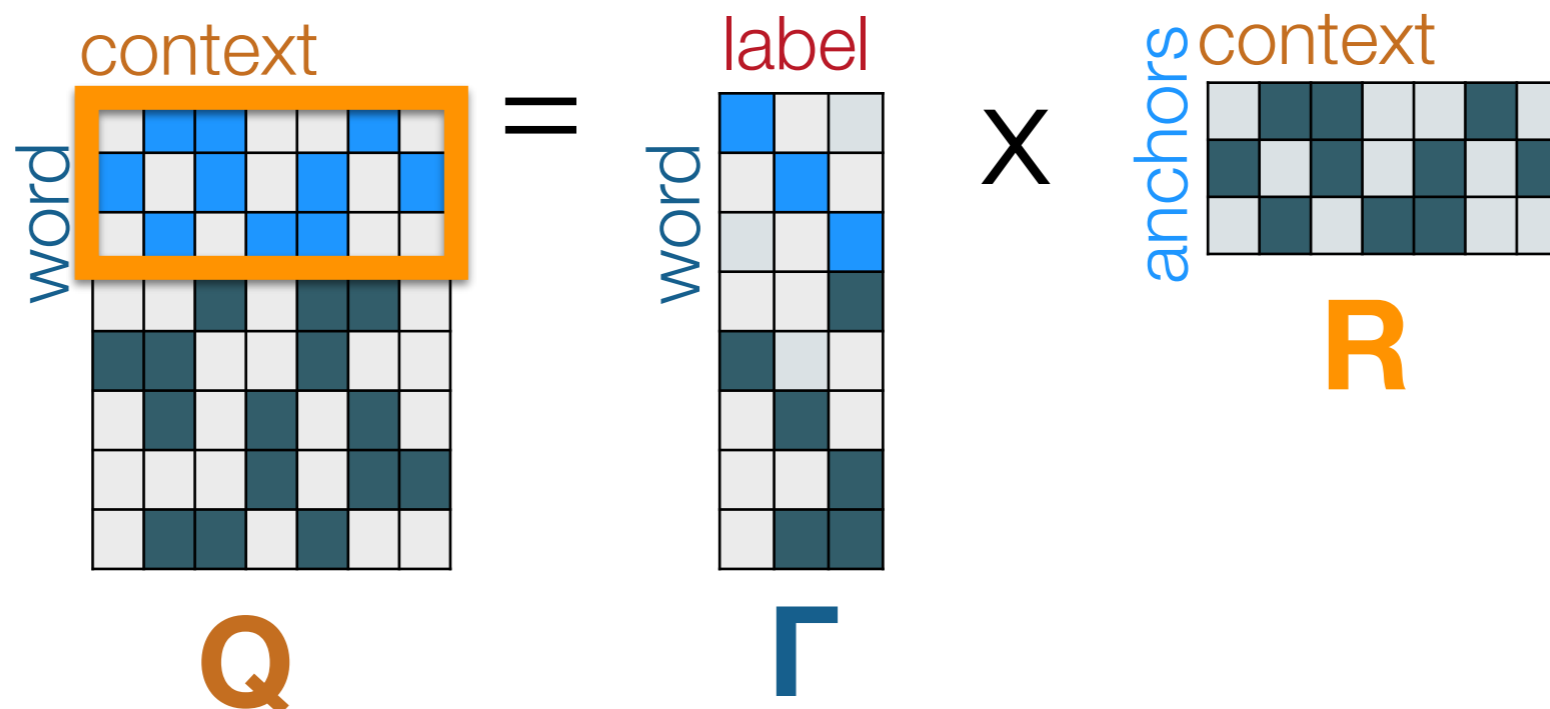


1. Conditional Independence $\text{word} \perp \text{context} \mid \text{label}$

$$p(\text{context} \mid \text{word}) = \sum_{\text{labels}} p(\text{label} \mid \text{word}) p(\text{context} \mid \text{label})$$

2. Anchor Trick

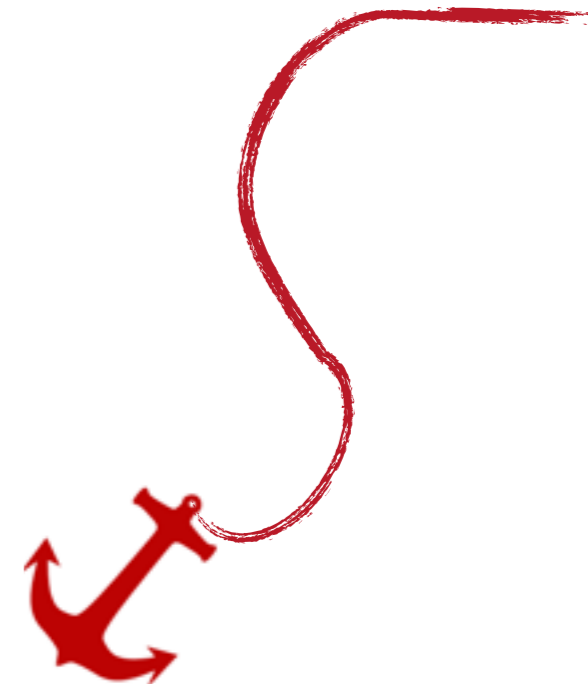
$$p(\text{context} \mid \text{word}) = \sum_{\text{labels}} p(\text{label} \mid \text{word}) p(\text{context} \mid \text{anchors})$$



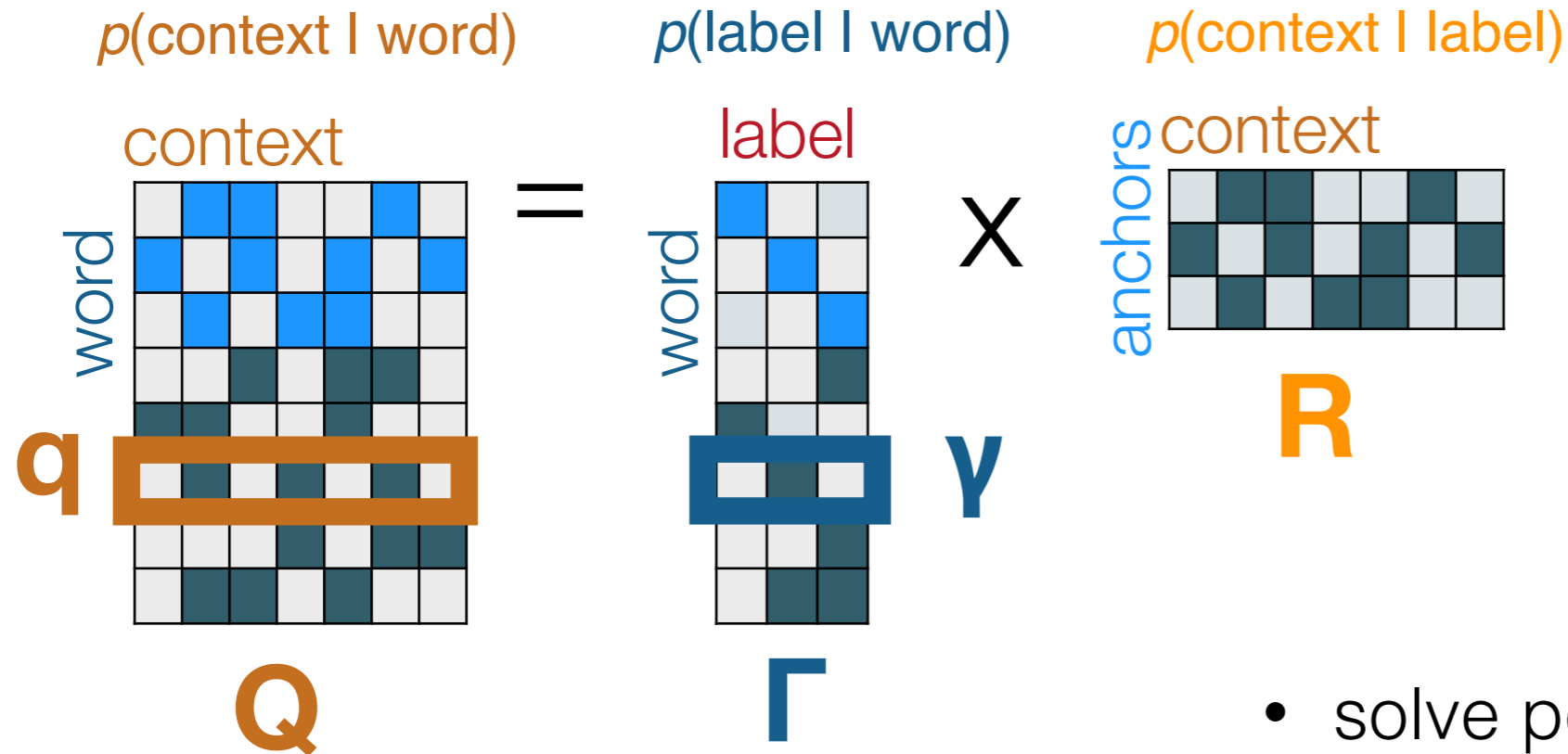
Learning sequence models via MoM

Outline

1. Learn HMM models via MoM
2. Solve a QP
3. Extend to feature-based model
4. Experiments



Method of Moments

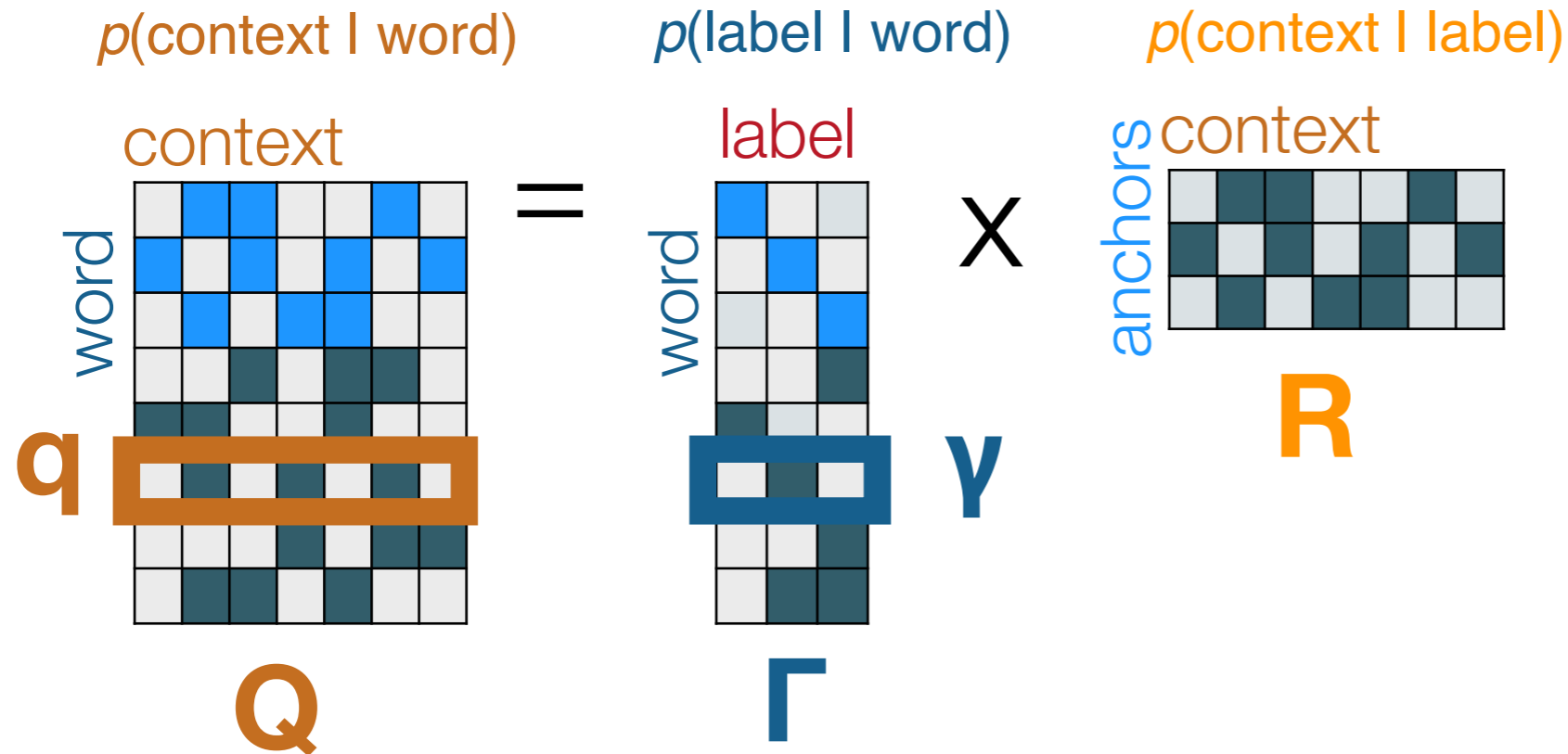


$$\boldsymbol{\gamma} = \operatorname{argmin} \|\mathbf{q} - \mathbf{R} \boldsymbol{\gamma}\|^2$$

$$0 \leq \boldsymbol{\gamma} \leq 1$$

$$\sum_{\text{labels}} \boldsymbol{\gamma} = 1$$

Method of Moments

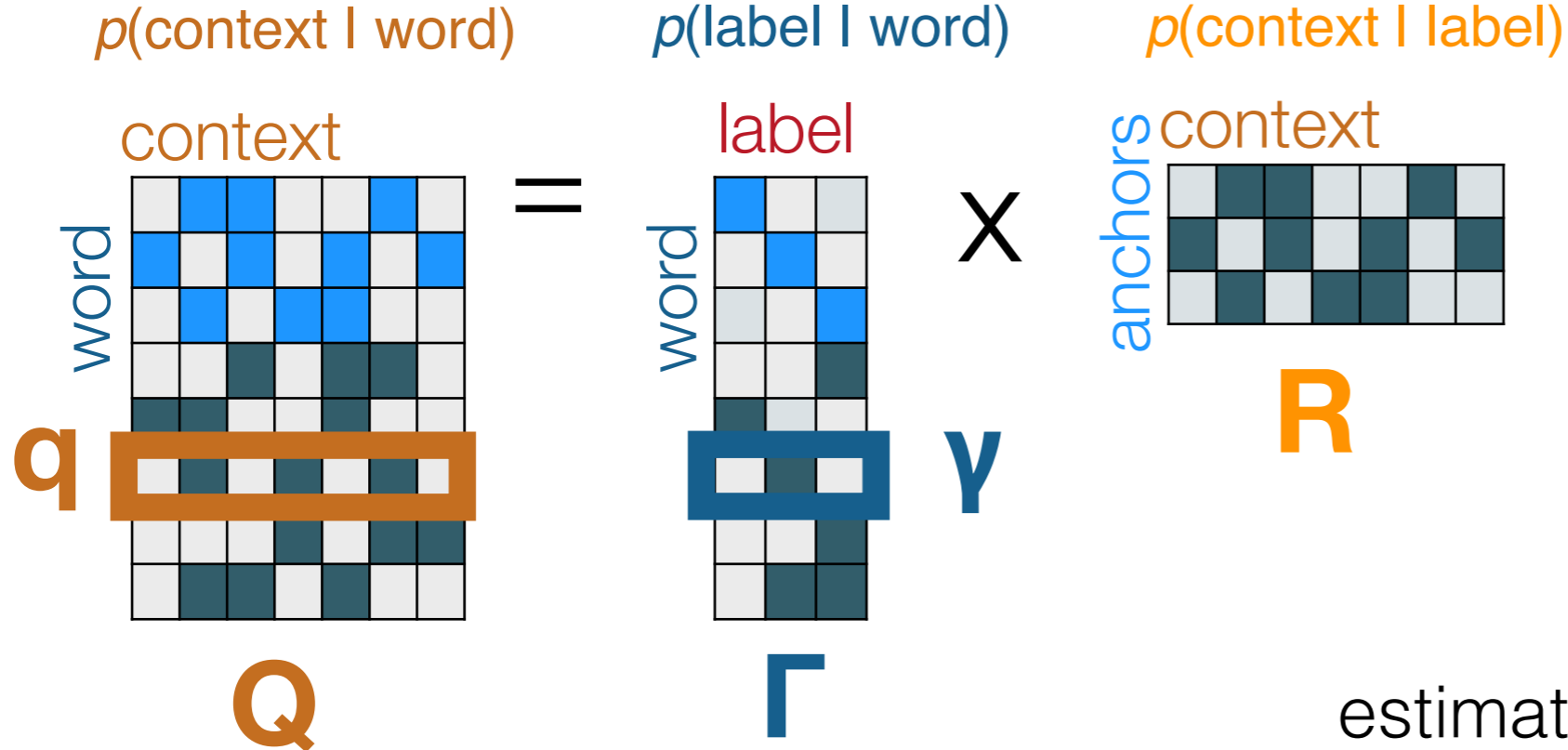


$$\mathbf{y} = \operatorname{argmin} \|\mathbf{q} - \mathbf{R} \mathbf{y}\|^2 + \lambda \|\mathbf{y}_{\text{sup}} - \mathbf{y}\|^2$$

$$0 \leq \mathbf{y} \leq 1$$

$$\sum_{\text{labels}} \mathbf{y} = 1$$

Method of Moments



estimated from labeled data

$$y = \operatorname{argmin} \| q - R y \|^2 + \lambda \| y_{\text{sup}} - y \|^2$$

$$0 \leq y \leq 1$$

$$\sum_{\text{labels}} y = 1$$

estimated from unlabeled data

Learn parameters ?

$p(\text{label} \mid \text{word})$

γ

coefficients



Observation Matrix

Bayes' Rule

$$p(\text{word} \mid \text{label}) = \gamma \frac{p(\text{word})}{p(\text{label})}$$

$$p(\text{label}) = \sum_{\text{words}} \gamma p(\text{word})$$

Learn parameters ?

Observation Matrix

Bayes' Rule

$$p(\text{word} \mid \text{label}) = \gamma \frac{p(\text{word})}{p(\text{label})}$$

Transition Matrix

- estimate from labeled data only

Learning sequence models via MoM

Outline

1. Learn HMM models via MoM
2. Relax the notion of *anchors*
3. Solve a QP
4. Experiments



Semi-supervised Twitter POS tagging



Twitter dataset

2.7 M unlabeled tweets
1000-100 labeled tweets \approx 200k words

12 Universal POS

x prt verb verb det adj noun
hehe its gonna b a good day

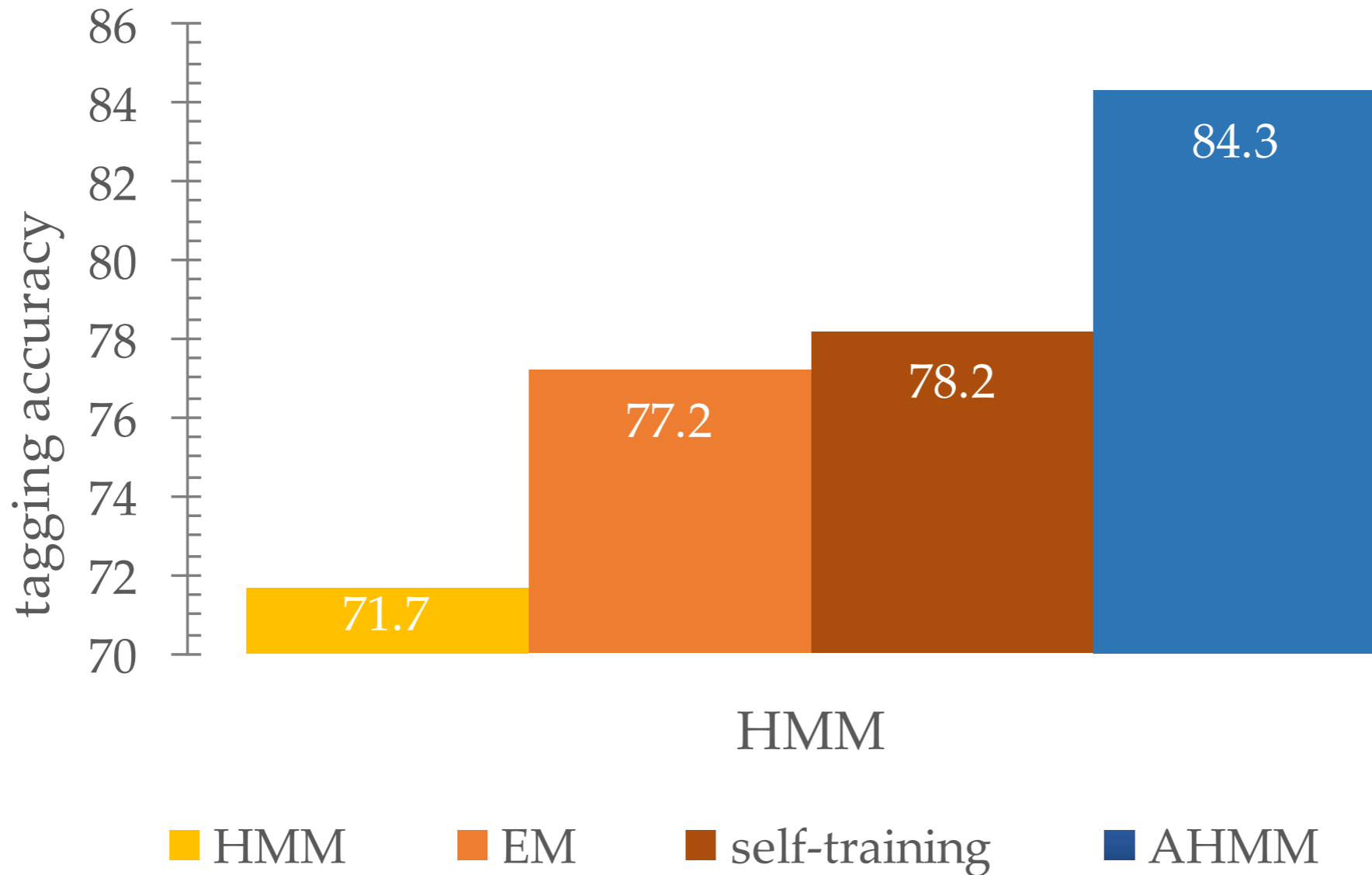
Slav Petrov et al., *A Universal Part-of-Speech Tagset*, 2011

Owoputi et al., Improved part-of-speech tagging for online conversational text with word clusters. 2013



Twitter POS tagging

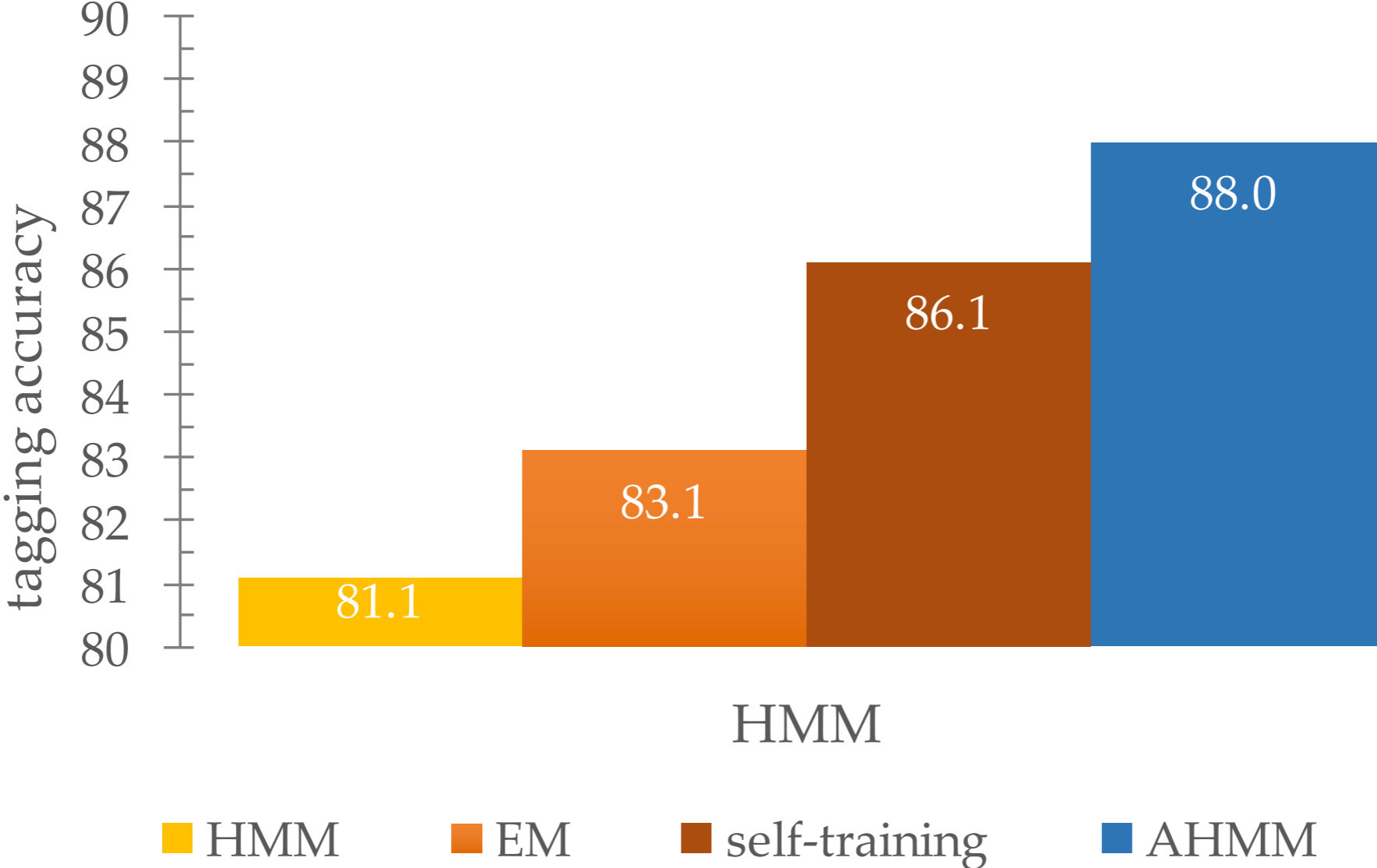
150 training labeled sequences





Twitter POS tagging

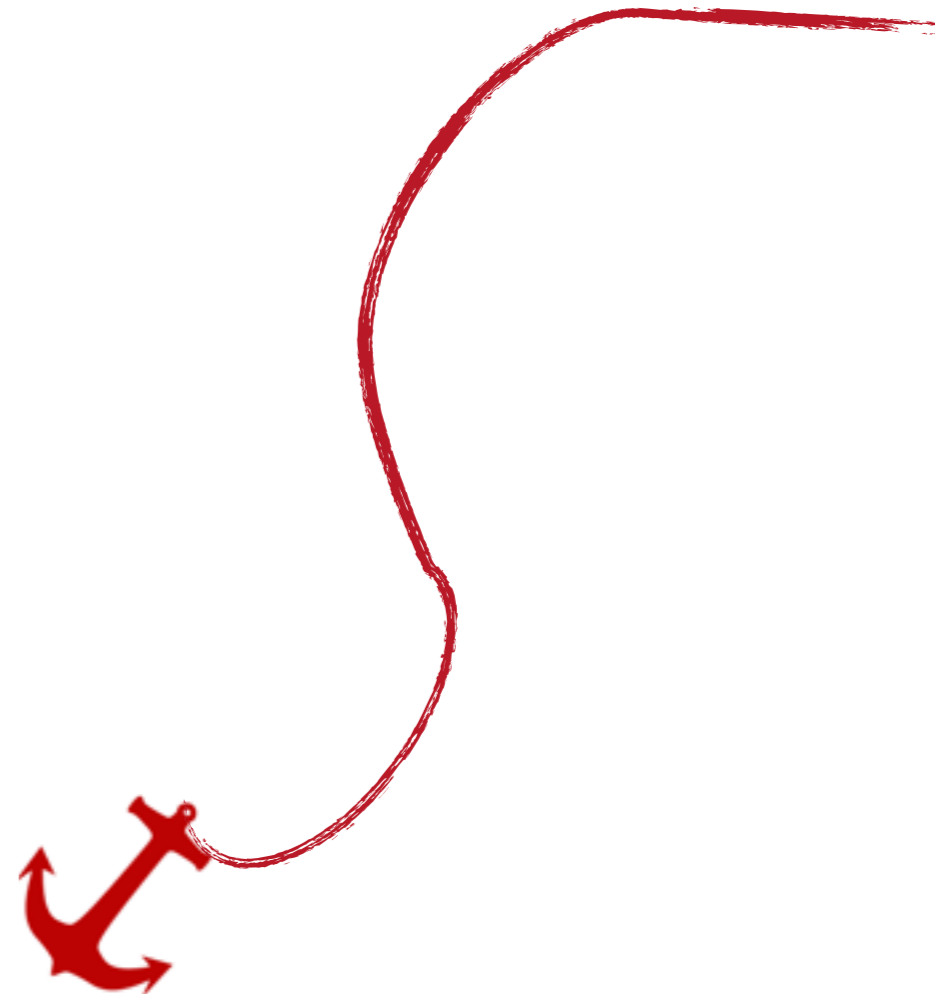
1000 training labeled sequences



Learning sequence models via MoM

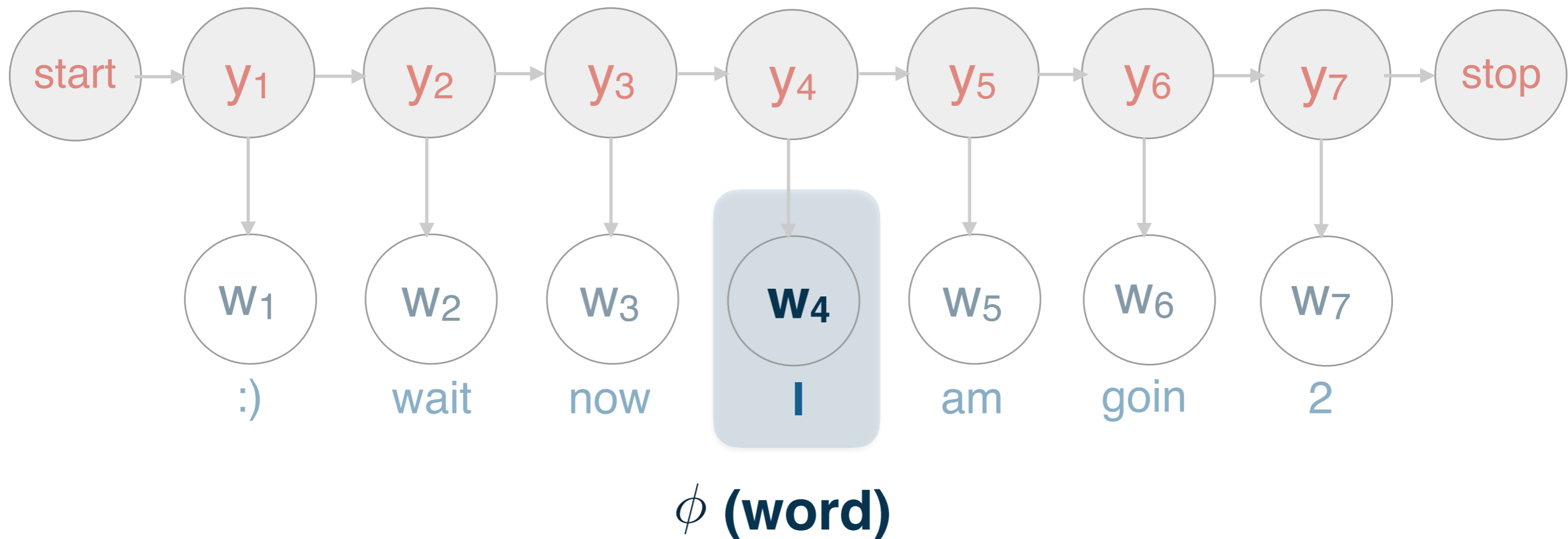
Outline

1. Learn HMM models via MoM
2. Relax the notion of *anchors*
3. Extend to feature HMM
4. Experiments



$\phi(\text{word})$

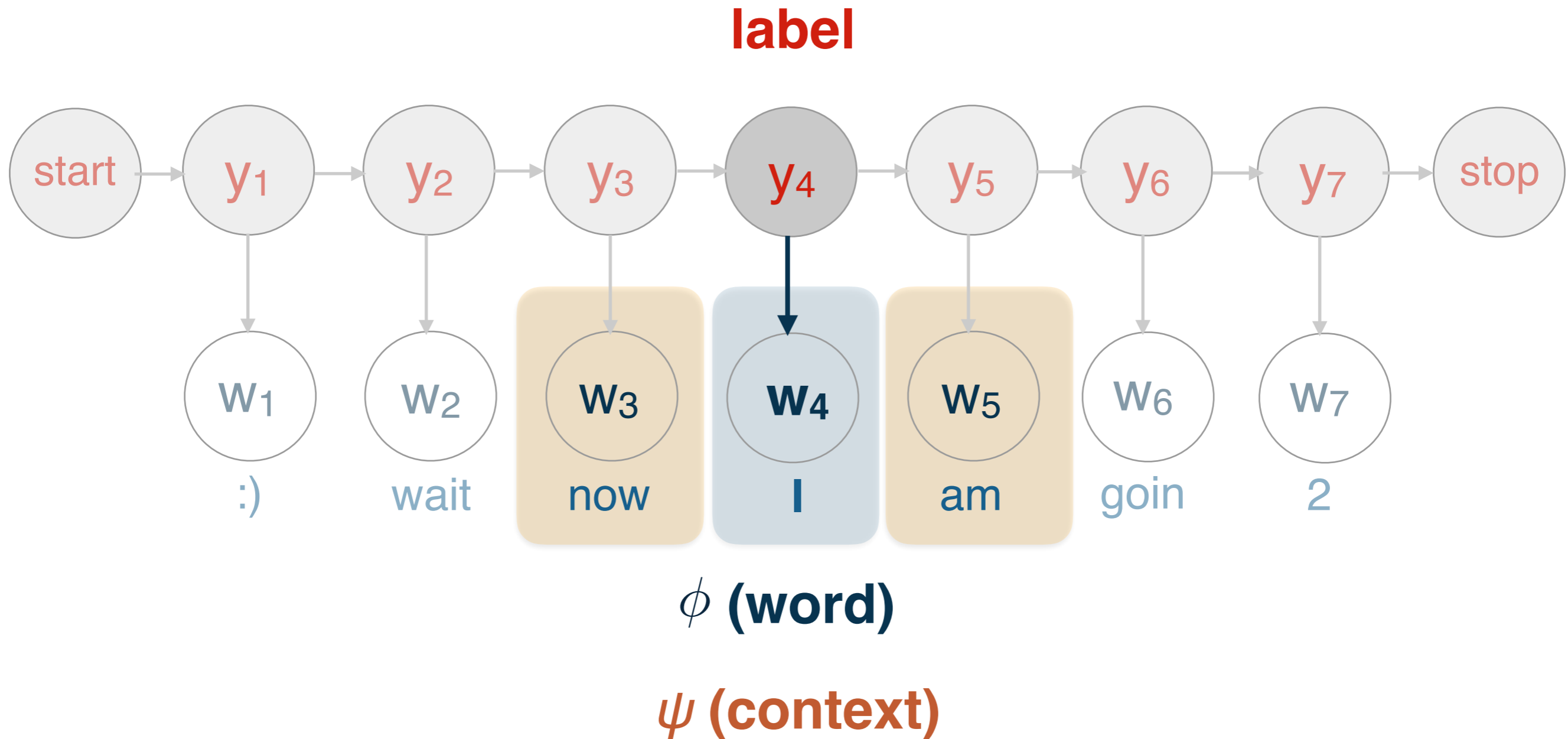
- is upper
- is title
- is digit
- is url
- starts #
- is emoticon



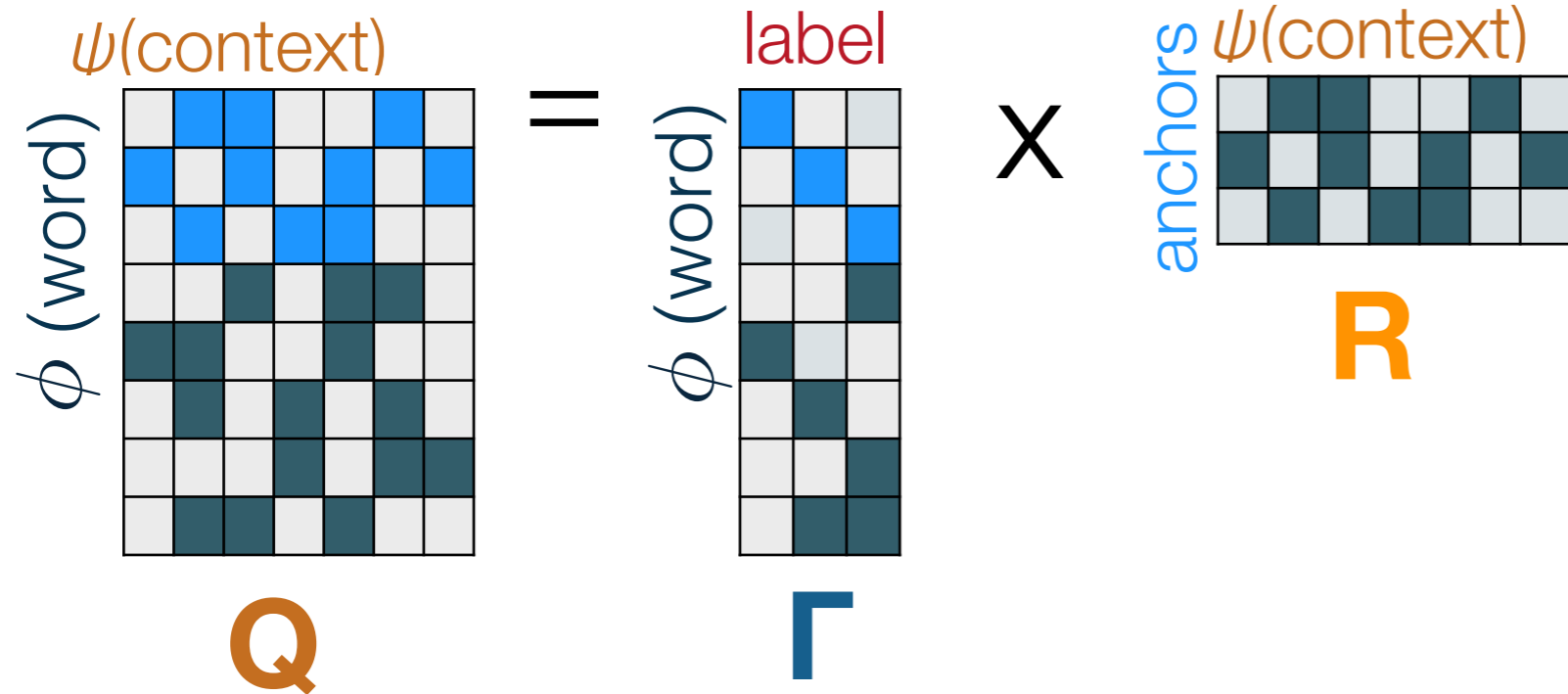
T. Berg-Kirkpatrick, *Painless unsupervised learning with features*, ACL 2010.

1. Conditional Independence

$$\text{word} \perp \text{context} \mid \text{label}$$



Log-linear model

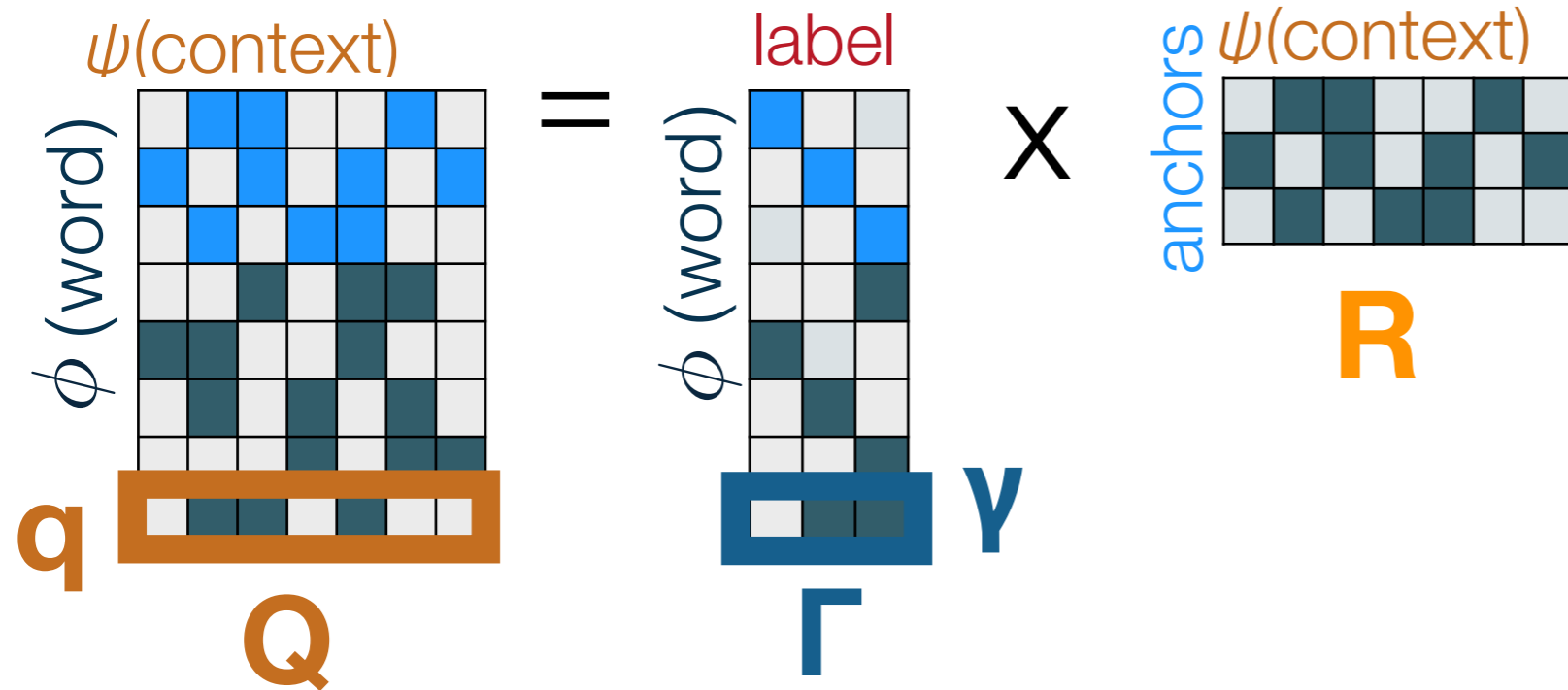


R $p(\text{context} \mid \text{label})$ \longrightarrow $E[\psi(\text{context}) \mid \text{label}]$

Q $p(\text{context} \mid \text{word})$ \longrightarrow $\frac{E[\psi(\text{context}) \times \Phi(\text{word})]}{E[\Phi(\text{word})]}$

Γ $p(\text{label} \mid \text{word})$ \longrightarrow $\frac{E[\Phi(\text{word}) \mid \text{label}] p(\text{label})}{E[\Phi(\text{word})]}$

Log-linear model



- solve per feature dimension Φ_j

$$\mathbf{y} = \operatorname{argmin} \|\mathbf{q} - \mathbf{R} \mathbf{y}\|^2 + \lambda \|\mathbf{y}_{\text{sup}} - \mathbf{y}\|^2$$

$$\sum_{\text{labels}} \mathbf{y} = 1$$

Learn parameters ?

$$\boldsymbol{\gamma} = \frac{E [\Phi(\text{word}) \mid \text{label}] p(\text{label})}{E [\Phi(\text{word})]}$$



mean parameters

$$\boldsymbol{\mu} = E[\Phi(\text{word}) \mid \text{label}] = \boldsymbol{\gamma} \frac{E[\Phi(\text{word})]}{p(\text{label})}$$

Learn parameters ?

mean parameters

canonical parameters

μ



θ_y

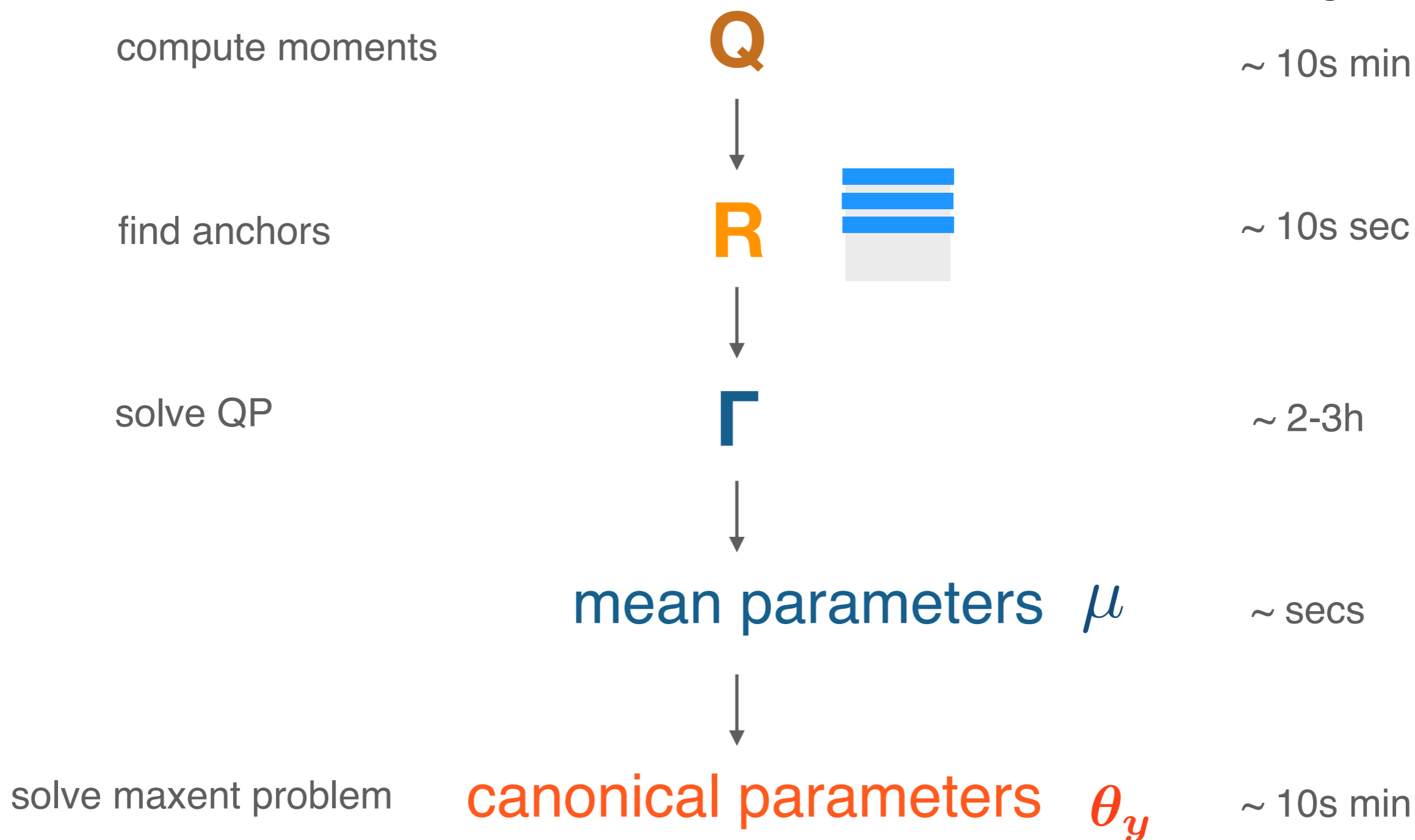
Fenchel-Legendre Duality

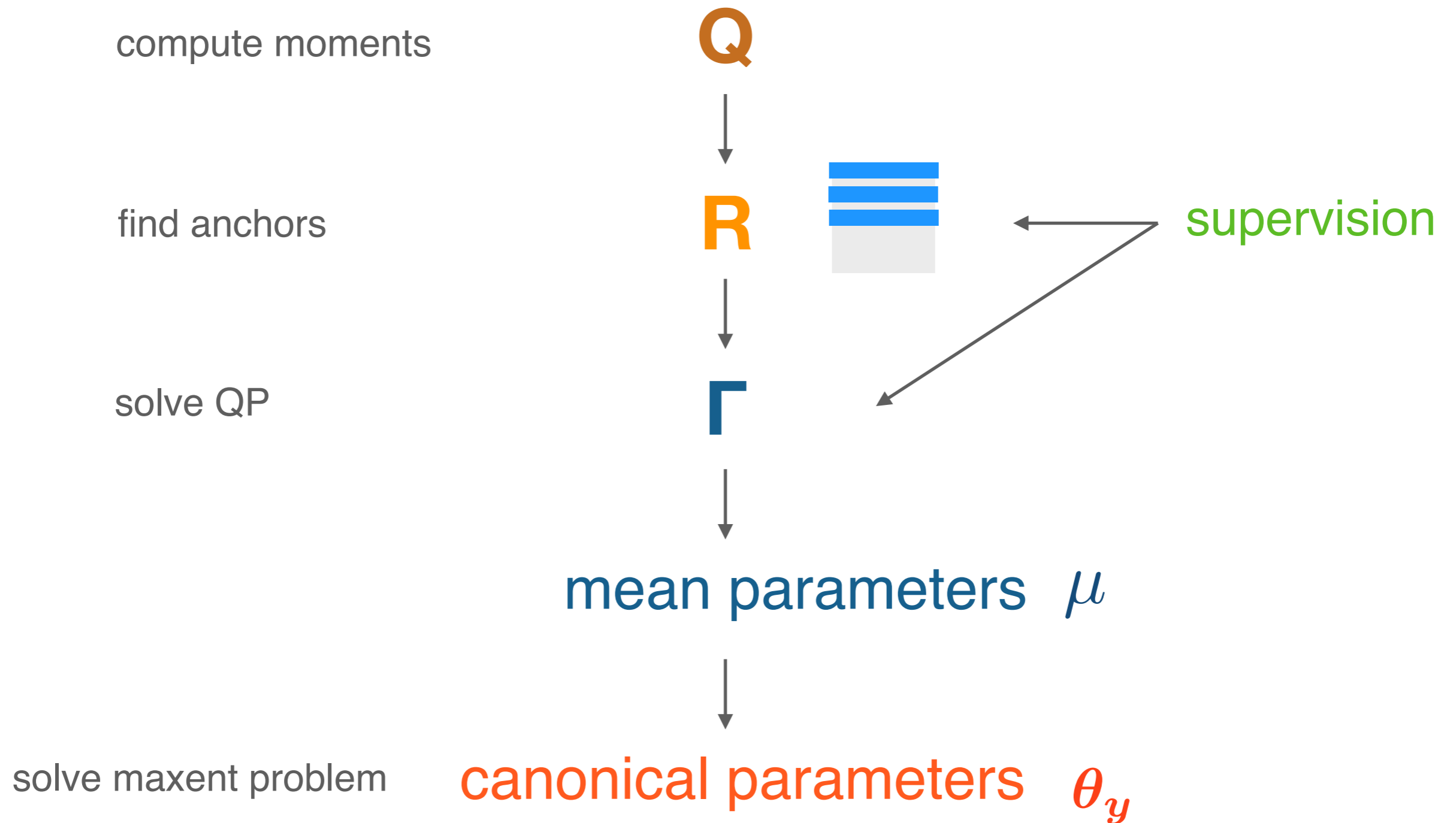
$$\theta_y^* = \operatorname{argmax}_{\theta_y} \theta_y^\top \mu_y - \log \mathcal{Z}_y$$

partition function

$$\mathcal{Z}_y = \sum_w \exp(\theta_y^\top t_w)$$

Algorithm





Learning sequence models via MoM

Outline

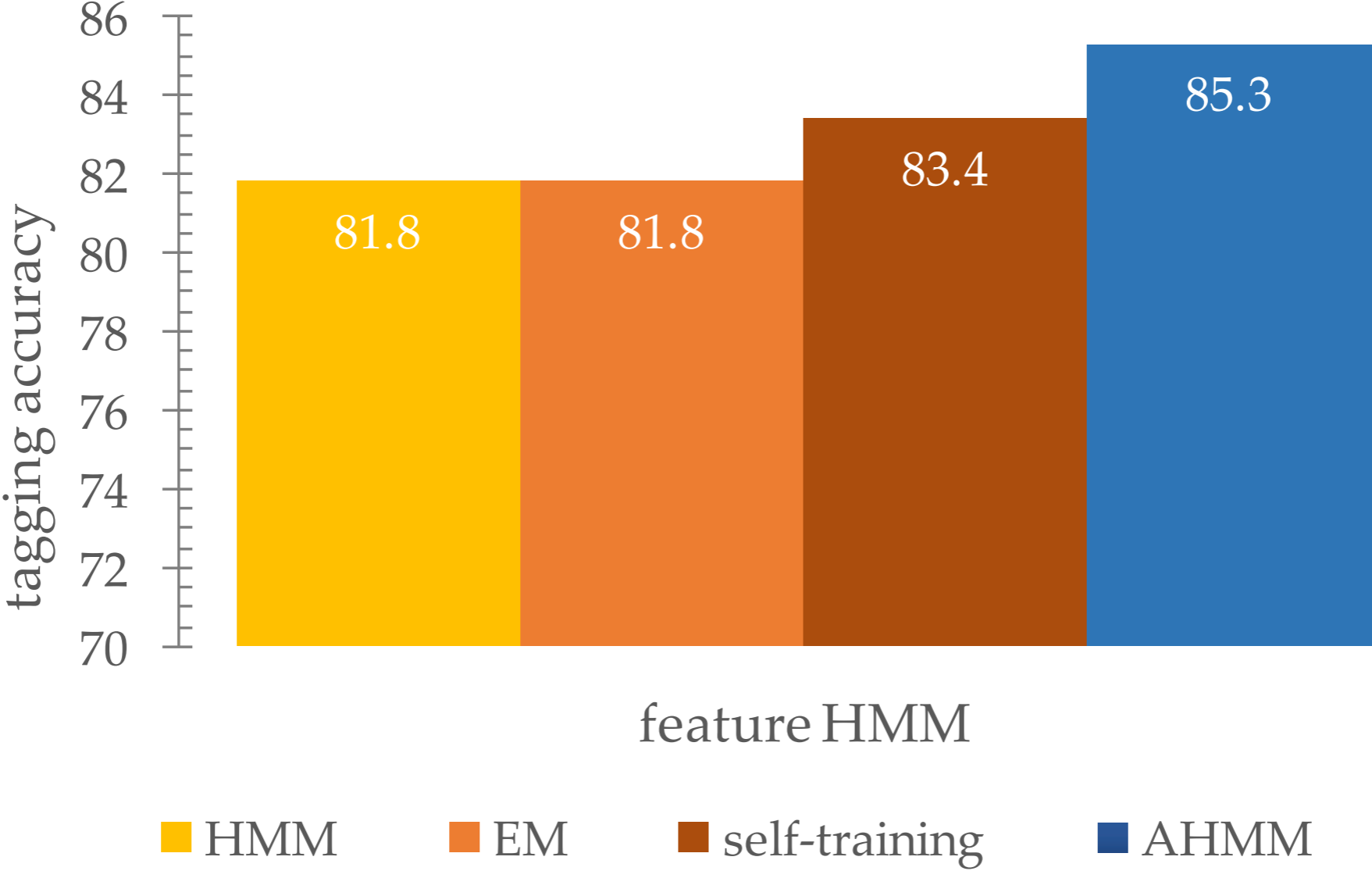
1. Learn HMM models via MoM
2. Relax the notion of *anchors*
3. Solve a QP
4. Experiments





Twitter POS tagging

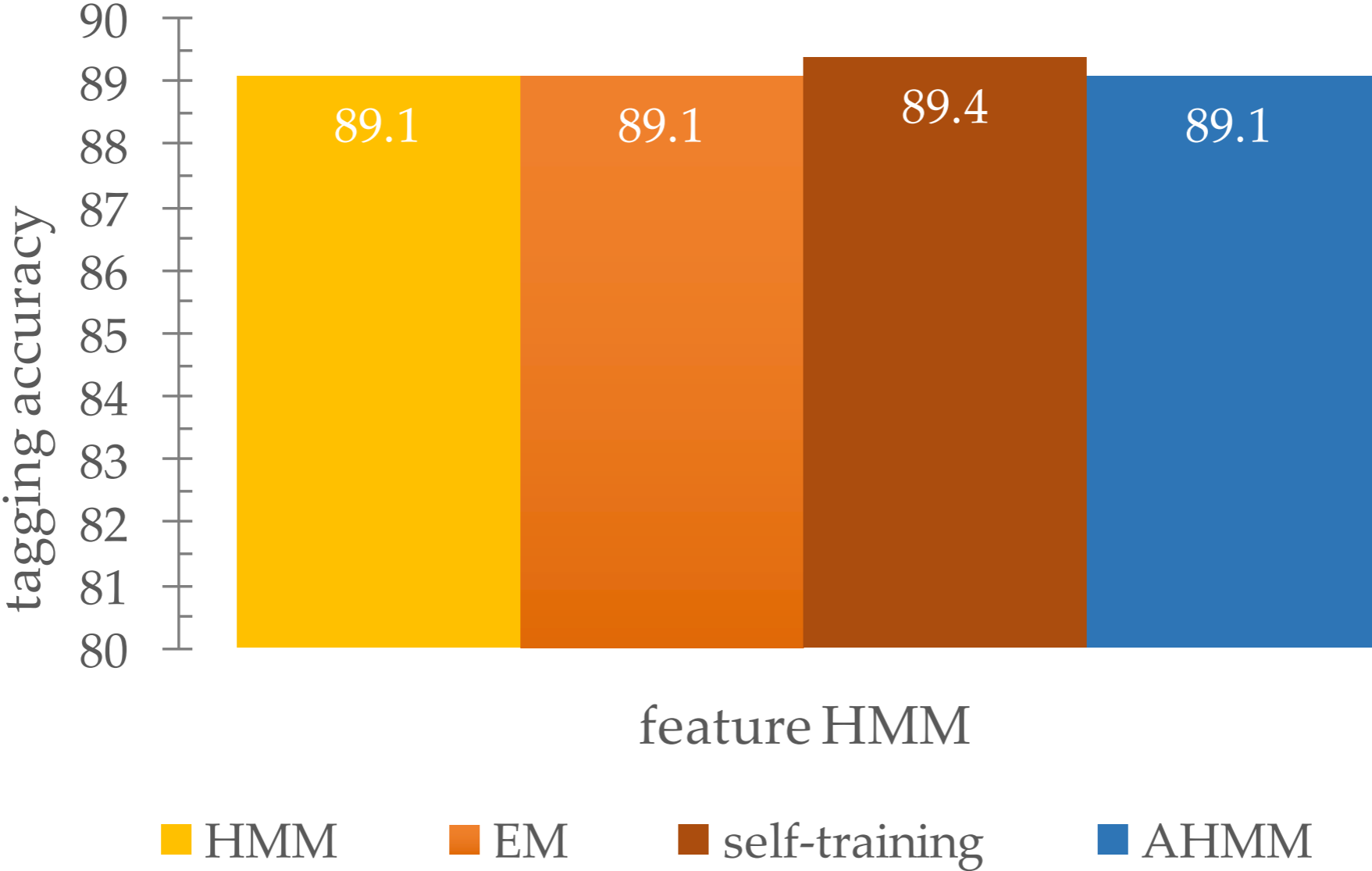
150 training labeled sequences





Twitter POS tagging

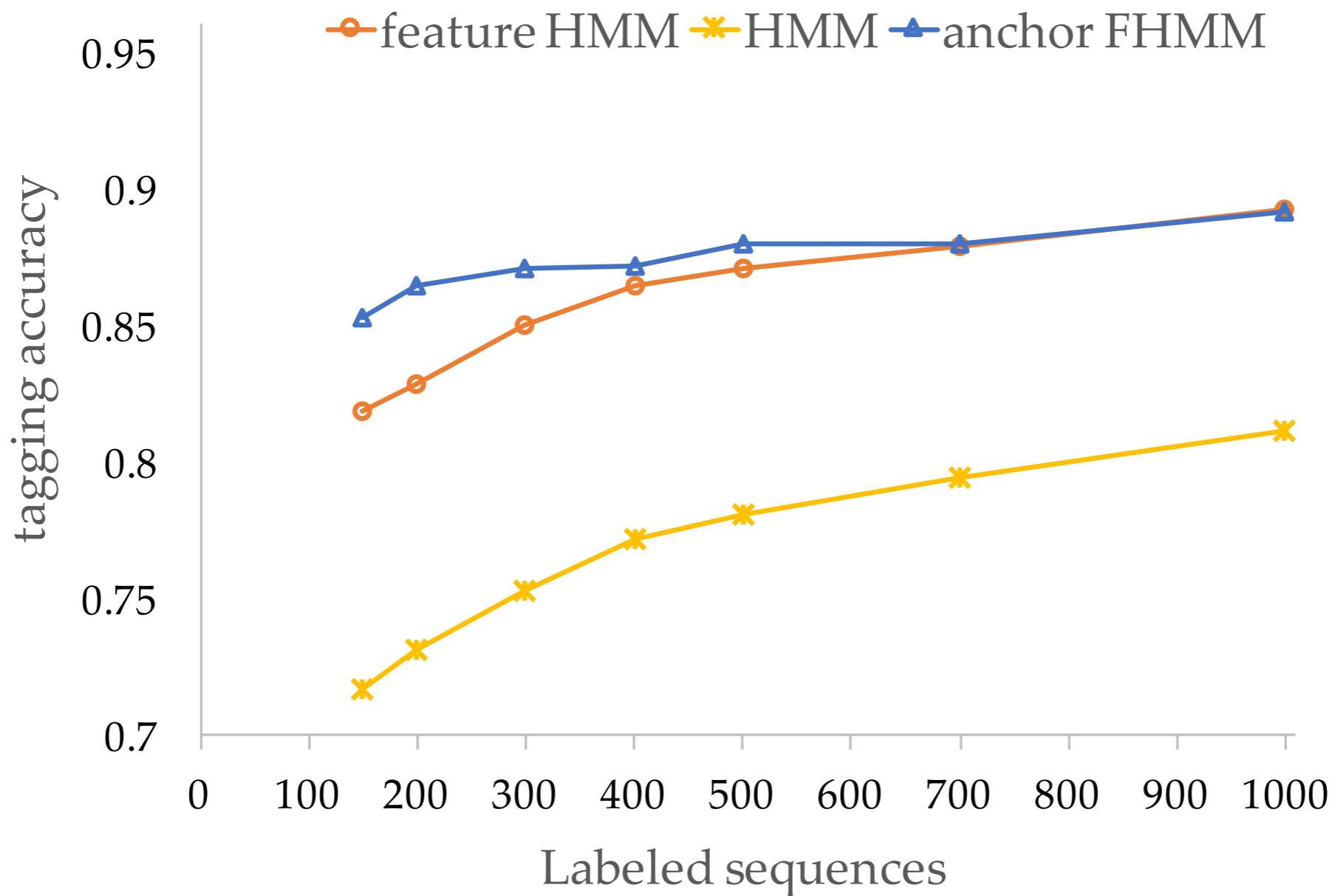
1000 training labeled sequences





Twitter POS tagging

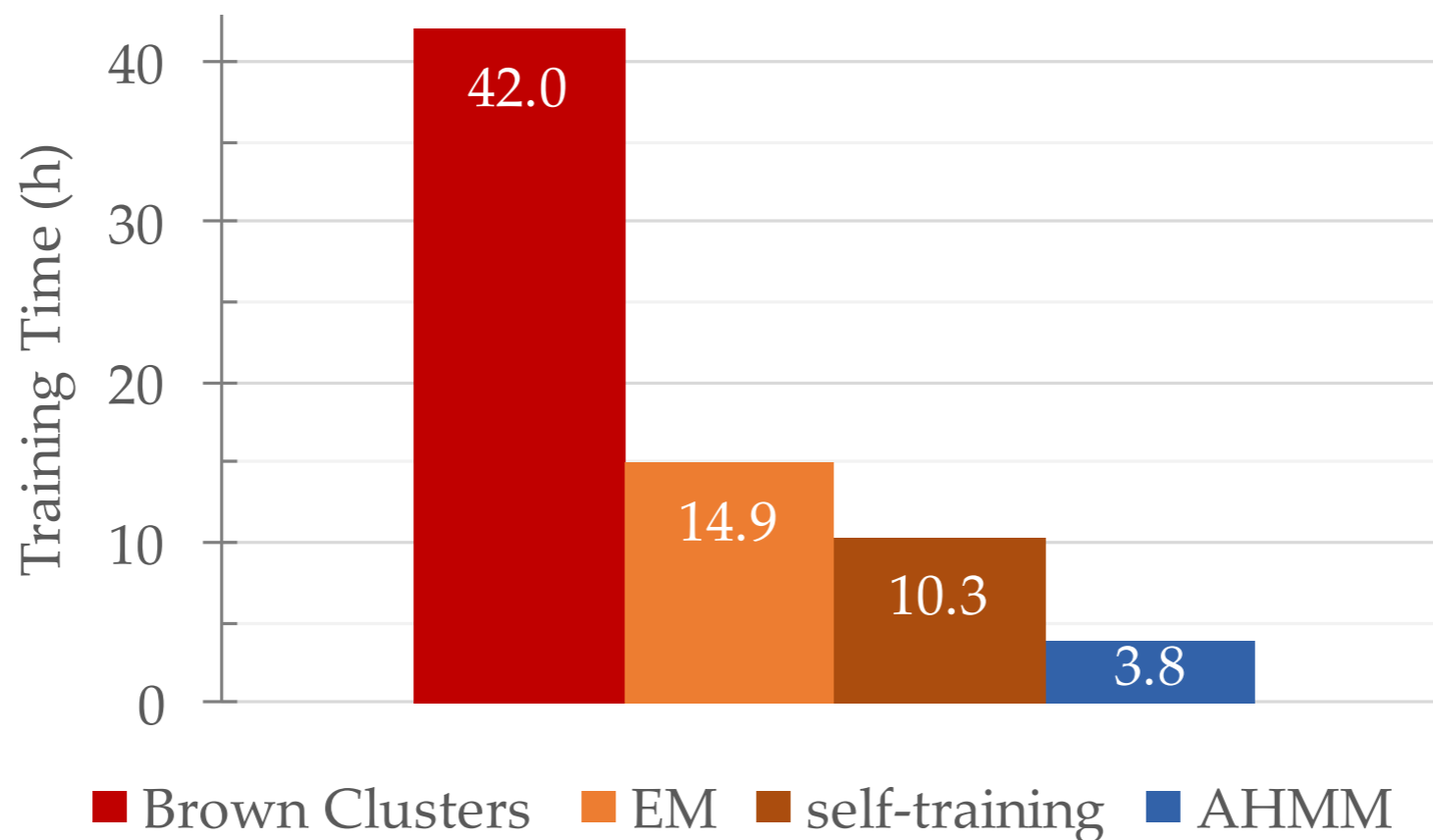
Tagging accuracy vs. labeled training size





Twitter POS tagging

1000 training sequences



Conclusions

- ✂ MoM algorithm for semi-supervised learning
- ✂ flexible method
(easy to add supervision)
- ✂ fast to train
(only one pass over the data)
- ✂ particularly good with little supervision

Thank you !

zmarinho@cmu.edu

Support for this research was provided by the Portuguese Science and Technology Foundation (FCT) and CMU Portugal Program, grant SFRH/BD/52015/2012. This work has also been partially supported by the European Union under H2020 project SUMMA, grant 688139, and by FCT, through contracts UID/EEA/50008/2013, through the LearnBig project (PTDC/EEISII/7092/2014), and the GoLocal project (grant CMUPERI/TIC/0046/2014).