

# **ANLP**

## **10 - CRFs (sequences, part II)**

David Schlangen

University of Potsdam, MSc Cognitive Systems  
Winter 2019 / 2020

# Industry Day 2019

## Idea

Every two years, the Applied CL group invites a number of local companies (related to language technology) to Uni Potsdam, so that

- students get an idea on current topics and research/development practices in industry
- contacts for internships, part-time or full-time jobs can be made

The 2019 event: Friday November 22nd, 10am to 4pm

- Venue: Campus Griebnitzsee, Building 6, Room S24 / S25

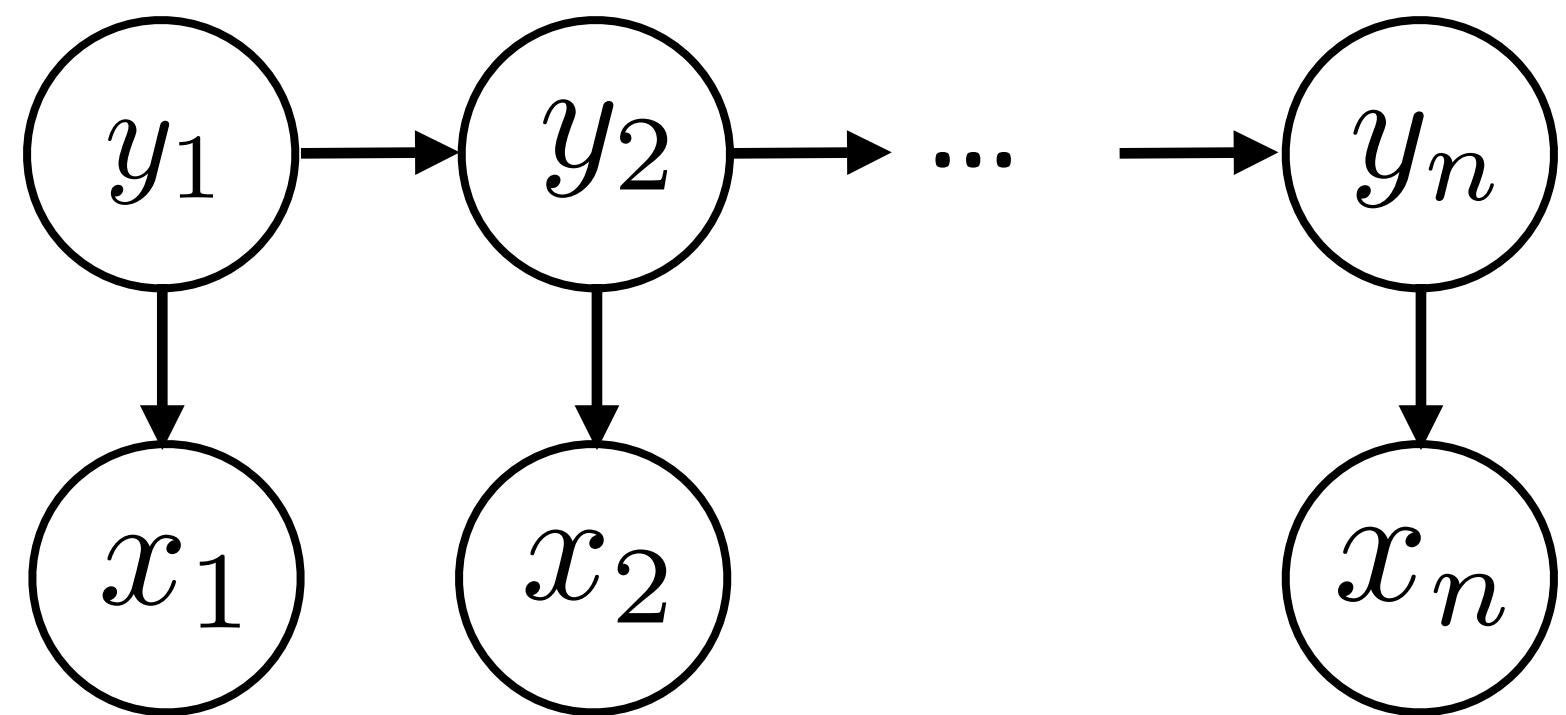
## Participating companies in 2019

- [acrolinx](#)
- [Carmeq](#)
- [DFKI](#)
- [ID Berlin](#)
- [parlamind](#)
- [rasa](#)
- [retresco](#)
- [zalando research](#)

# Recall: HMMs

---

- ▶ Input  $\mathbf{x} = (x_1, \dots, x_n)$       Output  $\mathbf{y} = (y_1, \dots, y_n)$



$$P(\mathbf{y}, \mathbf{x}) = P(y_1) \prod_{i=2}^n P(y_i|y_{i-1}) \prod_{i=1}^n P(x_i|y_i)$$

- ▶ Training: maximum likelihood estimation (with smoothing)

- ▶ Inference problem:  $\operatorname{argmax}_{\mathbf{y}} P(\mathbf{y}|\mathbf{x}) = \operatorname{argmax}_{\mathbf{y}} \frac{P(\mathbf{y}, \mathbf{x})}{P(\mathbf{x})}$

- ▶ Viterbi:  $\operatorname{score}_i(s) = \max_{y_{i-1}} P(s|y_{i-1}) P(x_i|s) \operatorname{score}_{i-1}(y_{i-1})$

# This Lecture

---

- ▶ CRFs: model (+features for NER), inference, learning
- ▶ Named entity recognition (NER)
- ▶ (if time) Beam search

# Named Entity Recognition

---

B-PER	I-PER	O	O	O	B-LOC	O	O	O	B-ORG	O	O
<i>Barack Obama</i>					<i>Hangzhou</i>				<i>G20</i>		

PERSON                                    LOC                                    ORG

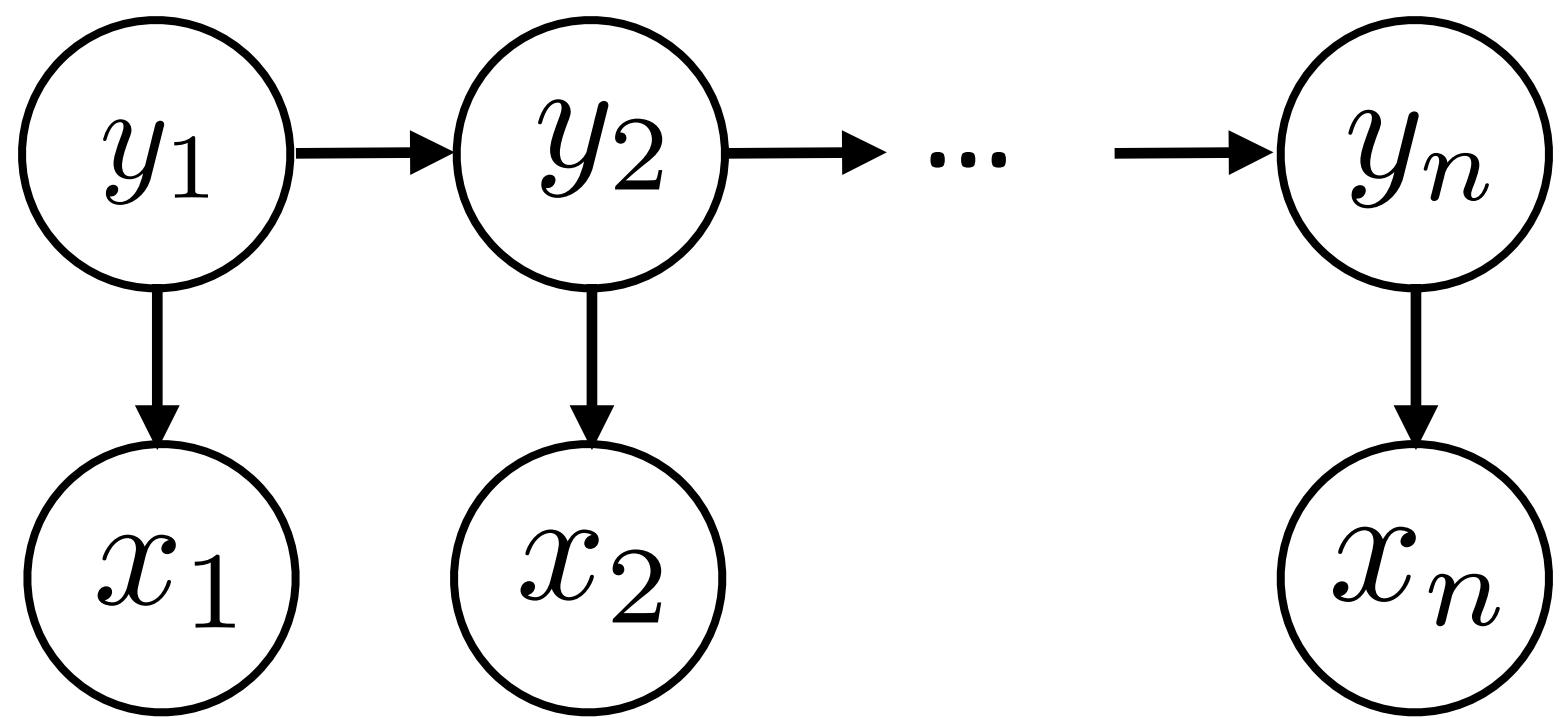
- ▶ BIO tagset: begin, inside, outside
- ▶ Sequence of tags — should we use an HMM?
- ▶ Why might an HMM not do so well here?
  - ▶ Lots of O's, so tags aren't as informative about context
  - ▶ Insufficient features/capacity with multinomials (especially for unks)

# CRFs

# Conditional Random Fields

---

- ▶ HMMs are expressible as Bayes nets (factor graphs)



- ▶ This reflects the following decomposition:

$$P(\mathbf{y}, \mathbf{x}) = P(y_1)P(x_1|y_1)P(y_2|y_1)P(x_2|y_2)\dots$$

- ▶ Locally normalized model:

- ▶ each factor is a probability distribution that normalizes

# Conditional Random Fields

---

- ▶ HMMs:  $P(\mathbf{y}, \mathbf{x}) = P(y_1)P(x_1|y_1)P(y_2|y_1)P(x_2|y_2) \dots$
- ▶ CRFs: discriminative models with the following globally-normalized form:

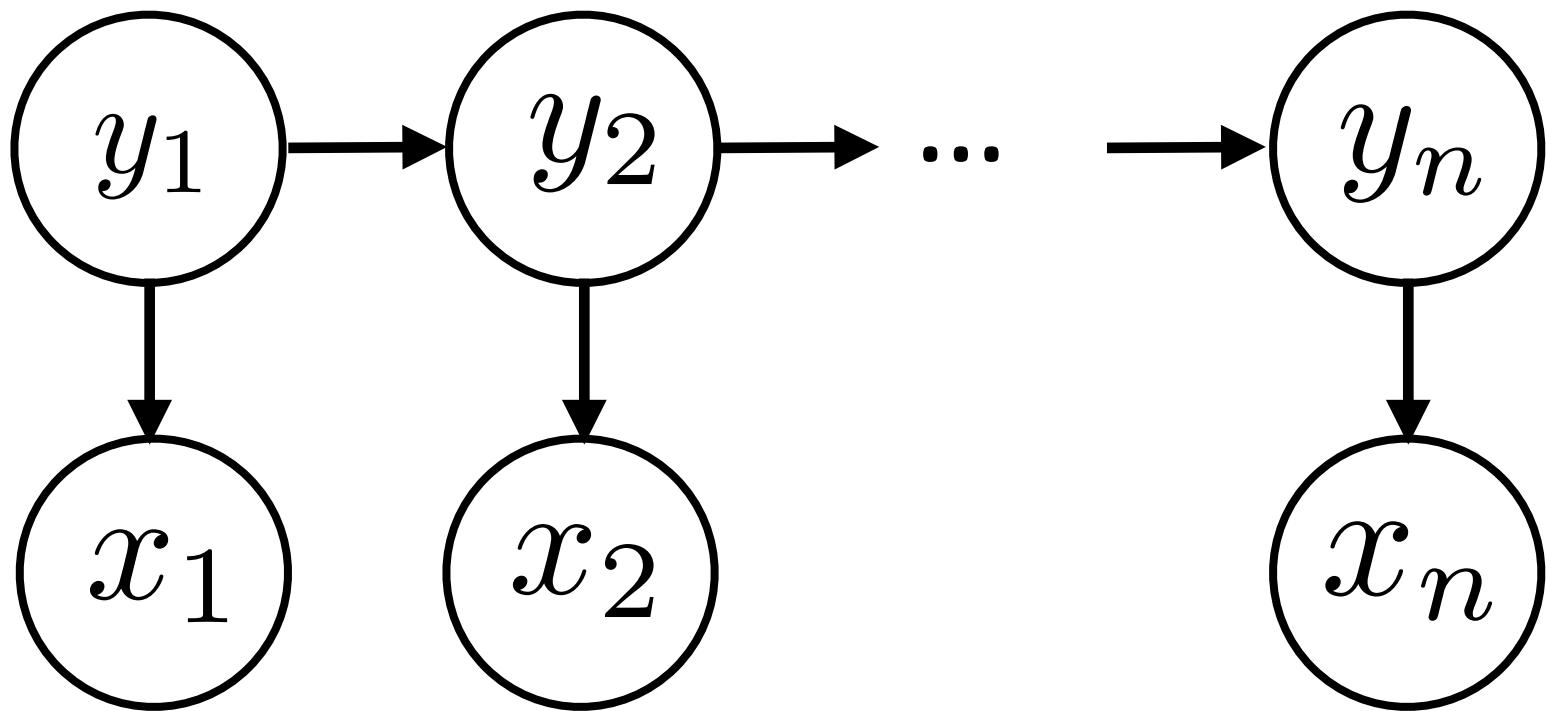
$$P(\mathbf{y}|\mathbf{x}) = \frac{1}{Z} \prod_k \exp(\phi_k(\mathbf{x}, \mathbf{y}))$$

↑  
normalizer                      any real-valued scoring function of its arguments

- ▶ Naive Bayes : logistic regression :: HMMs : CRFs  
local vs. global normalization  $\leftrightarrow$  generative vs. discriminative
- ▶ Locally normalized discriminative models do exist (MEMMs)
- ▶ How do we max over  $\mathbf{y}$ ? Intractable in general — can we fix this?

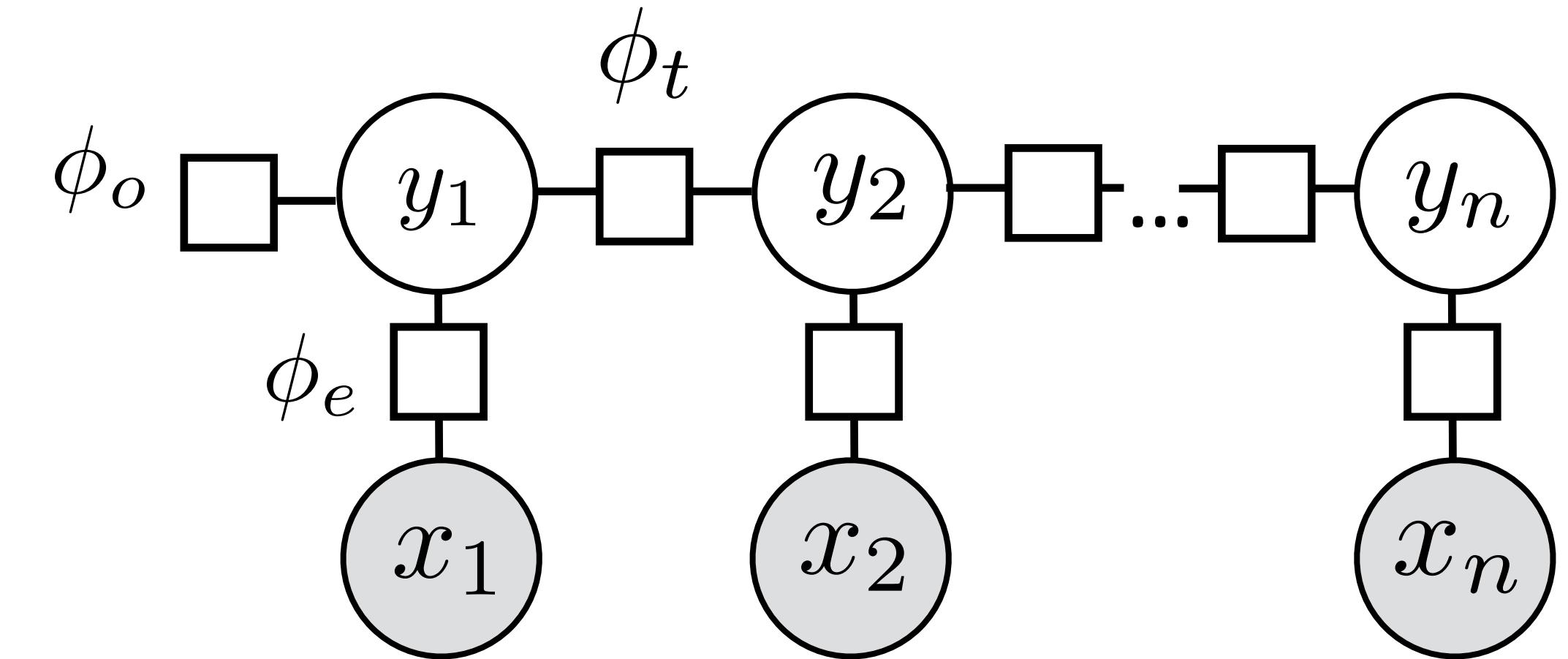
# Sequential CRFs

- ▶ HMMs:  $P(\mathbf{y}, \mathbf{x}) = P(y_1)P(x_1|y_1)P(y_2|y_1)P(x_2|y_2) \dots$



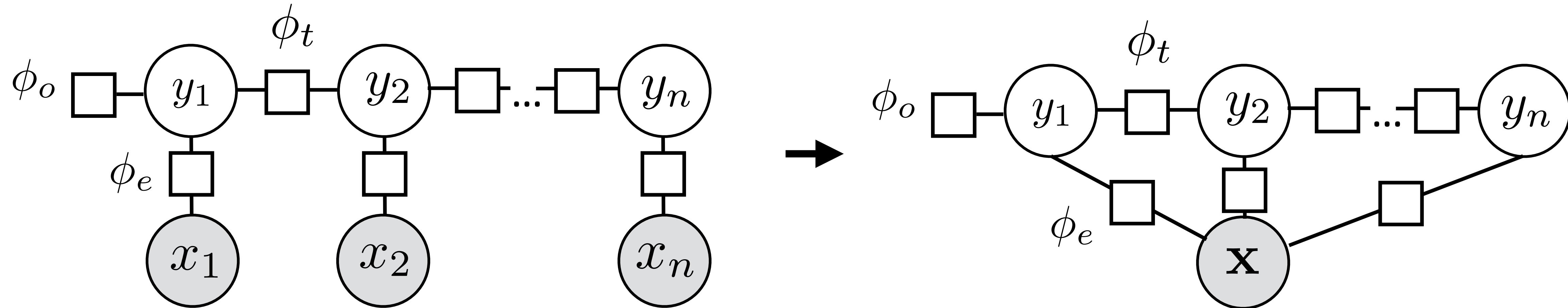
- ▶ CRFs:

$$P(\mathbf{y}|\mathbf{x}) \propto \prod_k \exp(\phi_k(\mathbf{x}, \mathbf{y}))$$



$$P(\mathbf{y}|\mathbf{x}) \propto \exp(\phi_o(y_1)) \prod_{i=2}^n \exp(\phi_t(y_{i-1}, y_i)) \prod_{i=1}^n \exp(\phi_e(x_i, y_i))$$

# Sequential CRFs

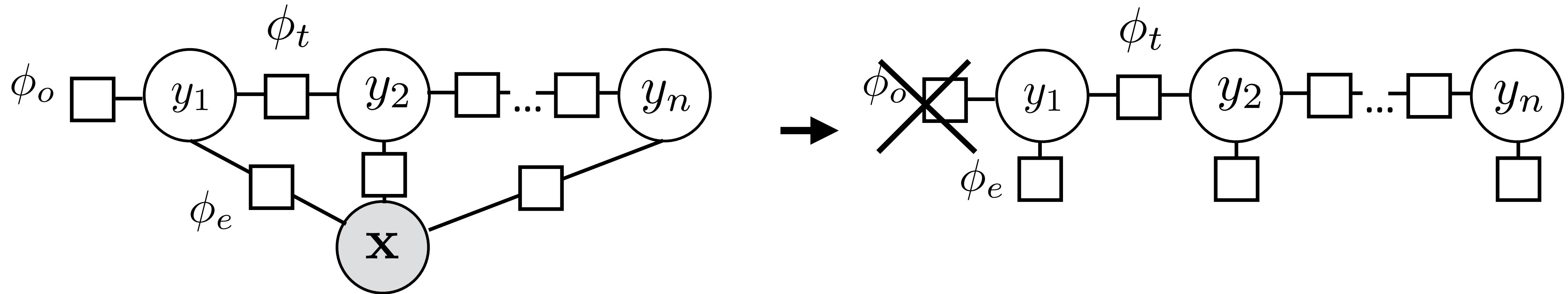


$$P(\mathbf{y}|\mathbf{x}) \propto \exp(\phi_o(y_1)) \prod_{i=2}^n \exp(\phi_t(y_{i-1}, y_i)) \prod_{i=1}^n \exp(\cancel{\phi_e(x_i, y_i)})$$

$$\prod_{i=1}^n \exp(\phi_e(y_i, i, \mathbf{x}))$$

- We condition on  $\mathbf{x}$ , so every factor can depend on all of  $\mathbf{x}$  (including transitions, but we won't do this)
  - $\mathbf{y}$  can't depend arbitrarily on  $\mathbf{x}$  in a generative model
- token index — lets us look at current word

# Sequential CRFs



- ▶ Notation: omit  $\mathbf{x}$  from the factor graph entirely (implicit)
- ▶ Don't include initial distribution, can bake into other factors

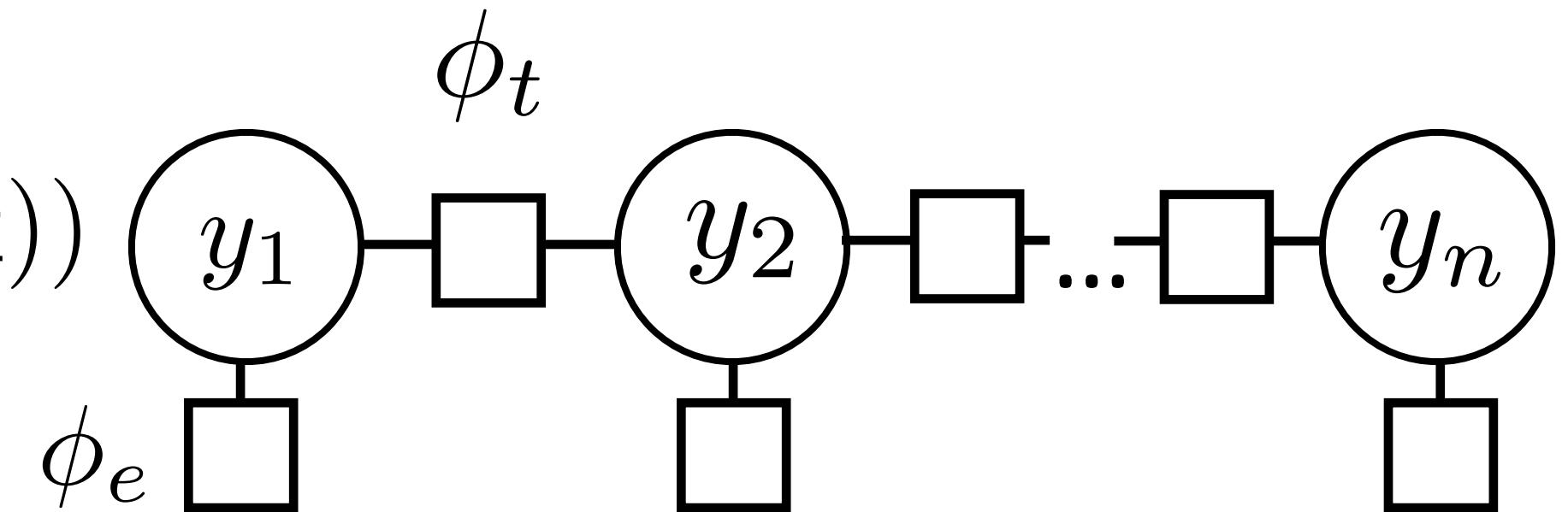
Sequential CRFs:

$$P(\mathbf{y}|\mathbf{x}) = \frac{1}{Z} \prod_{i=2}^n \exp(\phi_t(y_{i-1}, y_i)) \prod_{i=1}^n \exp(\phi_e(y_i, i, \mathbf{x}))$$

# Feature Functions

---

$$P(\mathbf{y}|\mathbf{x}) = \frac{1}{Z} \prod_{i=2}^n \exp(\phi_t(y_{i-1}, y_i)) \prod_{i=1}^n \exp(\phi_e(y_i, i, \mathbf{x}))$$



- This can be almost anything! Here we use linear functions of sparse features

$$\phi_e(y_i, i, \mathbf{x}) = w^\top f_e(y_i, i, \mathbf{x}) \quad \phi_t(y_{i-1}, y_i) = w^\top f_t(y_{i-1}, y_i)$$

$$P(\mathbf{y}|\mathbf{x}) \propto \exp w^\top \left[ \sum_{i=2}^n f_t(y_{i-1}, y_i) + \sum_{i=1}^n f_e(y_i, i, \mathbf{x}) \right]$$

- Looks like our single weight vector multiclass logistic regression model

$$\begin{aligned}
f(w = \text{they can fish}, y = \text{N V V}) &= \sum_{m=1}^{M+1} f(w, y_m, y_{m-1}, m) \\
&= f(w, \text{N}, \diamond, 1) \\
&\quad + f(w, \text{V}, \text{N}, 2) \\
&\quad + f(w, \text{V}, \text{V}, 3) \\
&\quad + f(w, \blacklozenge, \text{V}, 4) \\
&= (w_m = \text{they}, y_m = \text{N}) + (y_m = \text{N}, y_{m-1} = \diamond) \\
&\quad + (w_m = \text{can}, y_m = \text{V}) + (y_m = \text{V}, y_{m-1} = \text{N}) \\
&\quad + (w_m = \text{fish}, y_m = \text{V}) + (y_m = \text{V}, y_{m-1} = \text{V}) \\
&\quad + (y_m = \blacklozenge, y_{m-1} = \text{V}).
\end{aligned}$$

# Basic Features for NER

---

$$P(\mathbf{y}|\mathbf{x}) \propto \exp w^\top \left[ \sum_{i=2}^n f_t(y_{i-1}, y_i) + \sum_{i=1}^n f_e(y_i, i, \mathbf{x}) \right]$$

O      B-LOC  
*Barack Obama will travel to **Hangzhou** today for the G20 meeting .*

Transitions:  $f_t(y_{i-1}, y_i) = \text{Ind}[y_{i-1} \& y_i] = \text{Ind}[O - B-LOC]$

Emissions:  $f_e(y_6, 6, \mathbf{x}) = \text{Ind}[B-LOC \& \text{Current word} = Hangzhou]$   
 $\text{Ind}[B-LOC \& \text{Prev word} = to]$

# Features for NER

---

- ▶ Word features (can use in HMM)

- ▶ Capitalization

- ▶ Word shape

- ▶ Prefixes/suffixes

- ▶ Lexical indicators

- ▶ Context features (can't use in HMM!)

- ▶ Words before/after

- ▶ Tags before/after

- ▶ Word clusters

- ▶ Gazetteers

*Leicestershire*

*Boston*

*Apple released a new version...*

*According to the New York Times...*

# CRFs Outline

---

► Model:  $P(\mathbf{y}|\mathbf{x}) = \frac{1}{Z} \prod_{i=2}^n \exp(\phi_t(y_{i-1}, y_i)) \prod_{i=1}^n \exp(\phi_e(y_i, i, \mathbf{x}))$

$$P(\mathbf{y}|\mathbf{x}) \propto \exp w^\top \left[ \sum_{i=2}^n f_t(y_{i-1}, y_i) + \sum_{i=1}^n f_e(y_i, i, \mathbf{x}) \right]$$

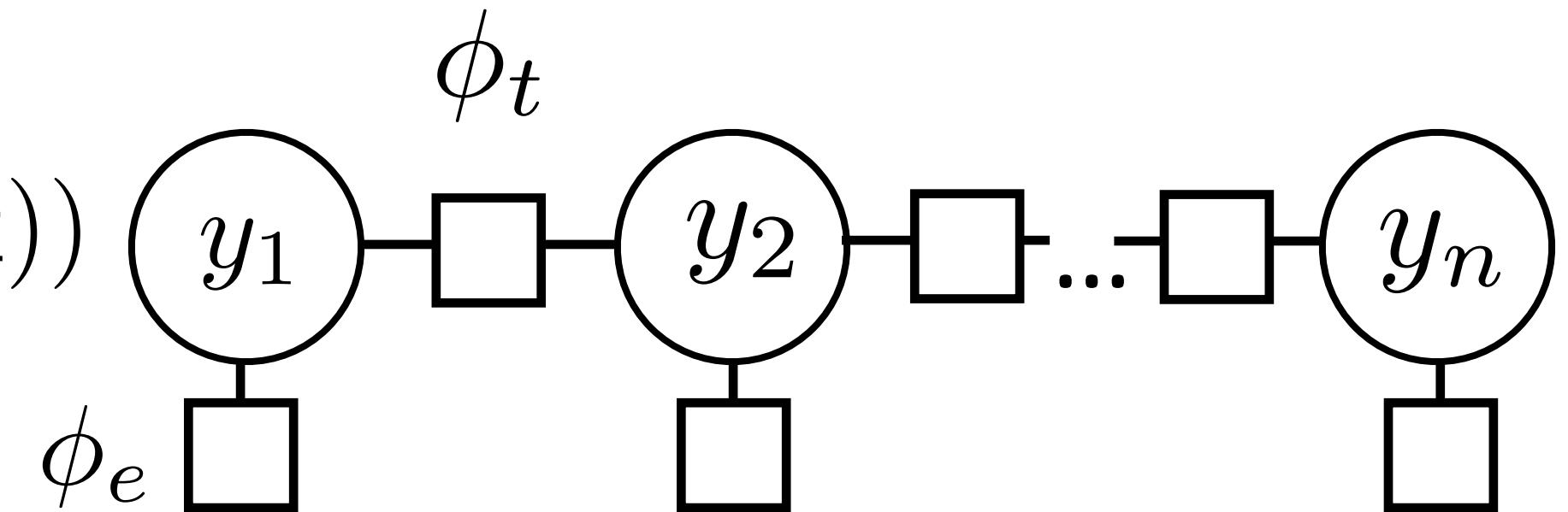
► Inference

► Learning

# Computing (arg)maxes

---

$$P(\mathbf{y}|\mathbf{x}) = \frac{1}{Z} \prod_{i=2}^n \exp(\phi_t(y_{i-1}, y_i)) \prod_{i=1}^n \exp(\phi_e(y_i, i, \mathbf{x}))$$



- $\text{argmax}_{\mathbf{y}} P(\mathbf{y}|\mathbf{x})$ : can use Viterbi exactly as in HMM case

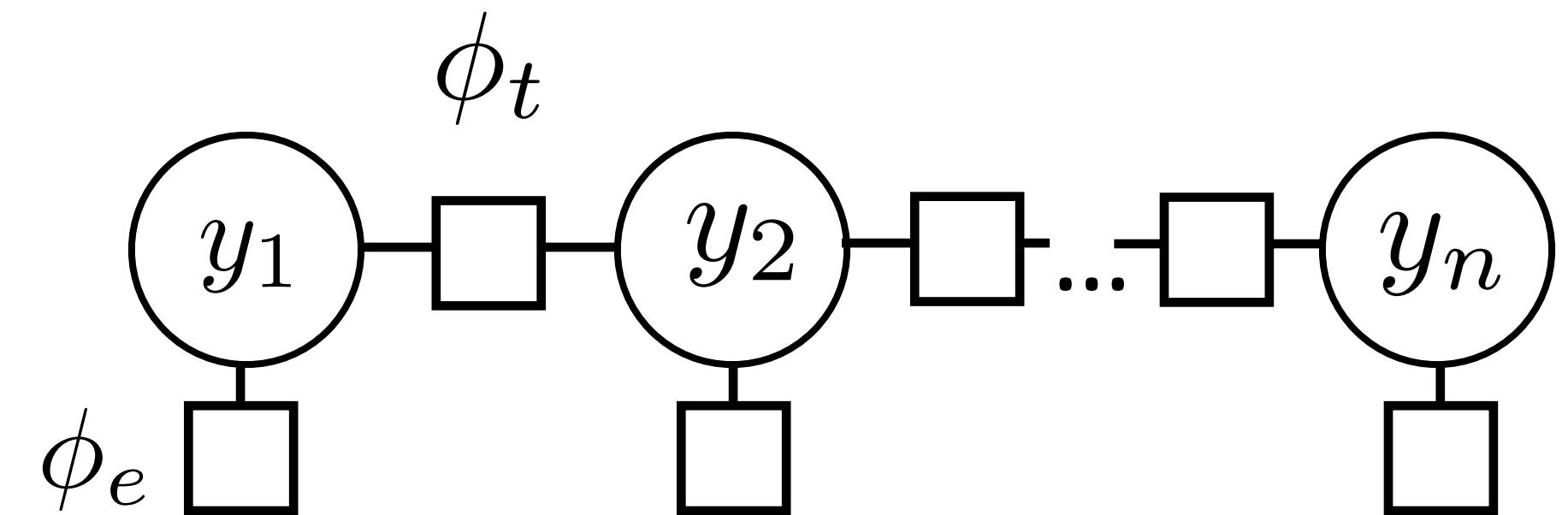
$$\begin{aligned} & \max_{y_1, \dots, y_n} e^{\phi_t(y_{n-1}, y_n)} e^{\phi_e(y_n, n, \mathbf{x})} \dots e^{\phi_e(y_2, 2, \mathbf{x})} e^{\phi_t(y_1, y_2)} e^{\phi_e(y_1, 1, \mathbf{x})} \\ &= \max_{y_2, \dots, y_n} e^{\phi_t(y_{n-1}, y_n)} e^{\phi_e(y_n, n, \mathbf{x})} \dots e^{\phi_e(y_2, 2, \mathbf{x})} \boxed{\max_{y_1} e^{\phi_t(y_1, y_2)}} \underbrace{e^{\phi_e(y_1, 1, \mathbf{x})}}_{\text{score}_1(y_1)} \\ &= \max_{y_3, \dots, y_n} e^{\phi_t(y_{n-1}, y_n)} e^{\phi_e(y_n, n, \mathbf{x})} \dots \max_{y_2} e^{\phi_t(y_2, y_3)} e^{\phi_e(y_2, 2, \mathbf{x})} \underbrace{\max_{y_1} e^{\phi_t(y_1, y_2)} \text{score}_1(y_1)}_{\text{score}_2(y_2)} \end{aligned}$$

- $\exp(\phi_t(y_{i-1}, y_i))$  and  $\exp(\phi_e(y_i, i, \mathbf{x}))$  play the role of the Ps now, same dynamic program

# Inference in General CRFs

---

- ▶ Can do inference in any tree-structured CRF



- ▶ Max-product algorithm: generalization of Viterbi to arbitrary tree-structured graphs (sum-product is generalization of forward-backward)

# CRFs Outline

---

- Model:  $P(\mathbf{y}|\mathbf{x}) = \frac{1}{Z} \prod_{i=2}^n \exp(\phi_t(y_{i-1}, y_i)) \prod_{i=1}^n \exp(\phi_e(y_i, i, \mathbf{x}))$   
$$P(\mathbf{y}|\mathbf{x}) \propto \exp w^\top \left[ \sum_{i=2}^n f_t(y_{i-1}, y_i) + \sum_{i=1}^n f_e(y_i, i, \mathbf{x}) \right]$$
- Inference: argmax  $P(\mathbf{y}|\mathbf{x})$  from Viterbi
- Learning

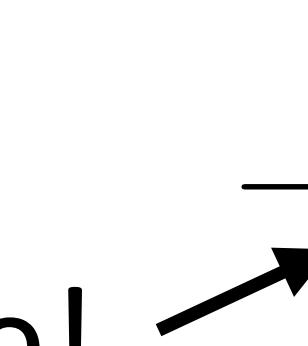
# Training CRFs

---

$$P(\mathbf{y}|\mathbf{x}) \propto \exp w^\top \left[ \sum_{i=2}^n f_t(y_{i-1}, y_i) + \sum_{i=1}^n f_e(y_i, i, \mathbf{x}) \right]$$

- ▶ Logistic regression:  $P(y|x) \propto \exp w^\top f(x, y)$
- ▶ Maximize  $\mathcal{L}(\mathbf{y}^*, \mathbf{x}) = \log P(\mathbf{y}^*|\mathbf{x})$
- ▶ Gradient is completely analogous to logistic regression:

$$\frac{\partial}{\partial w} \mathcal{L}(\mathbf{y}^*, \mathbf{x}) = \sum_{i=2}^n f_t(y_{i-1}^*, y_i^*) + \sum_{i=1}^n f_e(y_i^*, i, \mathbf{x})$$

intractable!   $-\mathbb{E}_{\mathbf{y}} \left[ \sum_{i=2}^n f_t(y_{i-1}, y_i) + \sum_{i=1}^n f_e(y_i, i, \mathbf{x}) \right]$

# Training CRFs

---

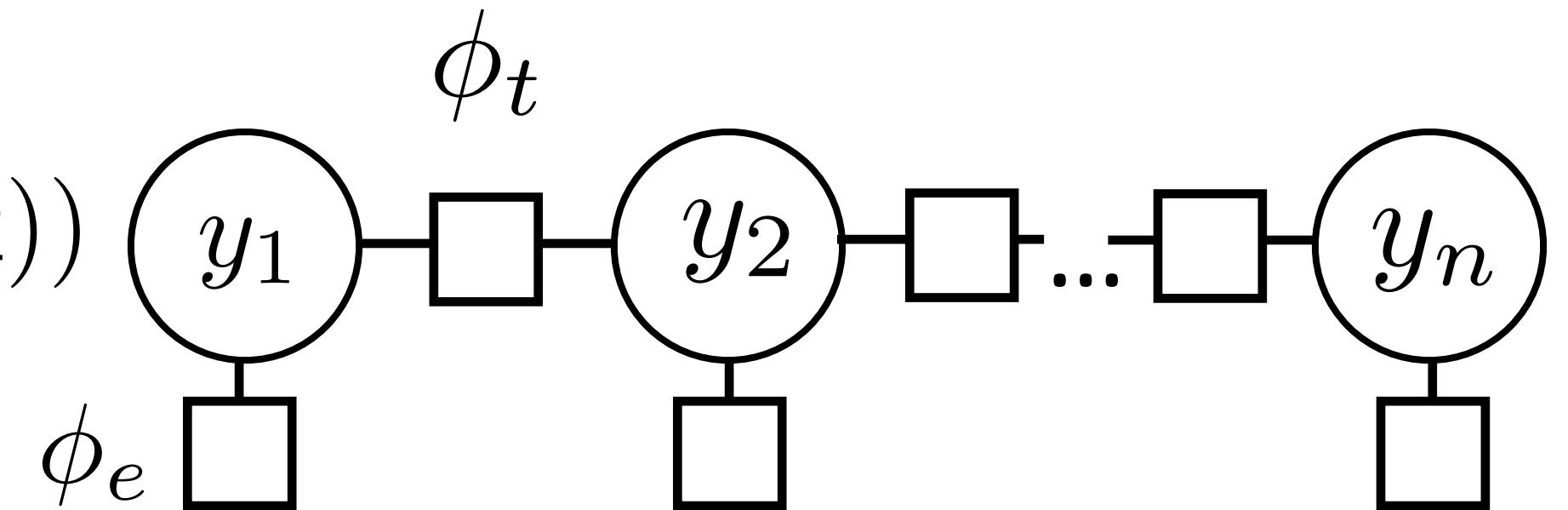
$$\begin{aligned}\frac{\partial}{\partial w} \mathcal{L}(\mathbf{y}^*, \mathbf{x}) &= \sum_{i=2}^n f_t(y_{i-1}^*, y_i^*) + \sum_{i=1}^n f_e(y_i^*, i, \mathbf{x}) \\ &\quad - \mathbb{E}_{\mathbf{y}} \left[ \sum_{i=2}^n f_t(y_{i-1}, y_i) + \sum_{i=1}^n f_e(y_i, i, \mathbf{x}) \right]\end{aligned}$$

► Let's focus on emission feature expectation

$$\begin{aligned}\mathbb{E}_{\mathbf{y}} \left[ \sum_{i=1}^n f_e(y_i, i, \mathbf{x}) \right] &= \sum_{\mathbf{y} \in \mathcal{Y}} P(\mathbf{y} | \mathbf{x}) \left[ \sum_{i=1}^n f_e(y_i, i, \mathbf{x}) \right] = \sum_{i=1}^n \sum_{\mathbf{y} \in \mathcal{Y}} P(\mathbf{y} | \mathbf{x}) f_e(y_i, i, \mathbf{x}) \\ &= \sum_{i=1}^n \sum_s P(y_i = s | \mathbf{x}) f_e(s, i, \mathbf{x})\end{aligned}$$

# Computing Marginals

$$P(\mathbf{y}|\mathbf{x}) = \frac{1}{Z} \prod_{i=2}^n \exp(\phi_t(y_{i-1}, y_i)) \prod_{i=1}^n \exp(\phi_e(y_i, i, \mathbf{x}))$$



- ▶ Normalizing constant  $Z = \sum_{\mathbf{y}} \prod_{i=2}^n \exp(\phi_t(y_{i-1}, y_i)) \prod_{i=1}^n \exp(\phi_e(y_i, i, \mathbf{x}))$
- ▶ Analogous to  $P(\mathbf{x})$  for HMMs
- ▶ For both HMMs and CRFs:

$$P(y_i = s | \mathbf{x}) = \frac{\text{forward}_i(s) \text{backward}_i(s)}{\sum_{s'} \text{forward}_i(s') \text{backward}_i(s')}$$

$Z$  for CRFs,  $P(\mathbf{x})$

for HMMs

$\text{backward}_i(s)$  = probability / score for generating rest of observation seq., from here

# Training CRFs

---

- ▶ For emission features:

$$\frac{\partial}{\partial w} \mathcal{L}(\mathbf{y}^*, \mathbf{x}) = \sum_{i=1}^n f_e(y_i^*, i, \mathbf{x}) - \sum_{i=1}^n \sum_s P(y_i = s | \mathbf{x}) f_e(s, i, \mathbf{x})$$

gold features – expected features under model

- ▶ Transition features: need to compute  $P(y_i = s_1, y_{i+1} = s_2 | \mathbf{x})$  using forward-backward as well
- ▶ ...but you can build a pretty good system without transition features

# CRFs Outline

---

► Model:  $P(\mathbf{y}|\mathbf{x}) = \frac{1}{Z} \prod_{i=2}^n \exp(\phi_t(y_{i-1}, y_i)) \prod_{i=1}^n \exp(\phi_e(y_i, i, \mathbf{x}))$

$$P(\mathbf{y}|\mathbf{x}) \propto \exp w^\top \left[ \sum_{i=2}^n f_t(y_{i-1}, y_i) + \sum_{i=1}^n f_e(y_i, i, \mathbf{x}) \right]$$

- Inference: argmax  $P(\mathbf{y}|\mathbf{x})$  from Viterbi
- Learning: run forward-backward to compute posterior probabilities; then

$$\frac{\partial}{\partial w} \mathcal{L}(\mathbf{y}^*, \mathbf{x}) = \sum_{i=1}^n f_e(y_i^*, i, \mathbf{x}) - \sum_{i=1}^n \sum_s P(y_i = s | \mathbf{x}) f_e(s, i, \mathbf{x})$$

NER

# NER

---

- ▶ CRF with lexical features can get around 85 F1 on this problem
- ▶ Other pieces of information that many systems capture
- ▶ World knowledge:

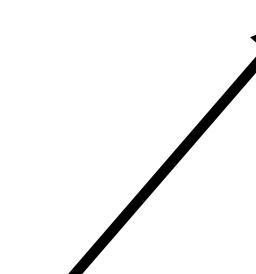
The delegation met the president at the airport, **Tanjug** said.

**Tanjug**

---

From Wikipedia, the free encyclopedia

**Tanjug** (/tʌnjʊg/) ([Serbian Cyrillic](#): Танјуг) is a Serbian state news agency based in [Belgrade](#).<sup>[2]</sup>



# Nonlocal Features

---

The news agency Tanjug reported on the outcome of the meeting.

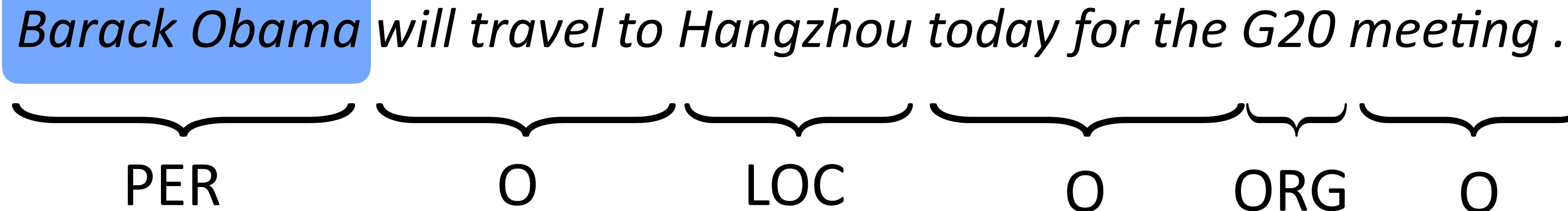
ORG?  
PER?

The delegation met the president at the airport, Tanjug said.

- More complex factor graph structures can let you capture this, or just decode sentences in order and use features on previous sentences

# Semi-Markov Models

---



- ▶ Chunk-level prediction rather than token-level BIO
- ▶  $y$  is a set of touching spans of the sentence
- ▶ Pros: features can look at whole span at once
- ▶ Cons: there's an extra factor of  $n$  in the dynamic programs

# Evaluating NER

---

B-PER I-PER O O O B-LOC O O O B-ORG O O

*Barack Obama will travel to Hangzhou today for the G20 meeting .*

PERSON

LOC

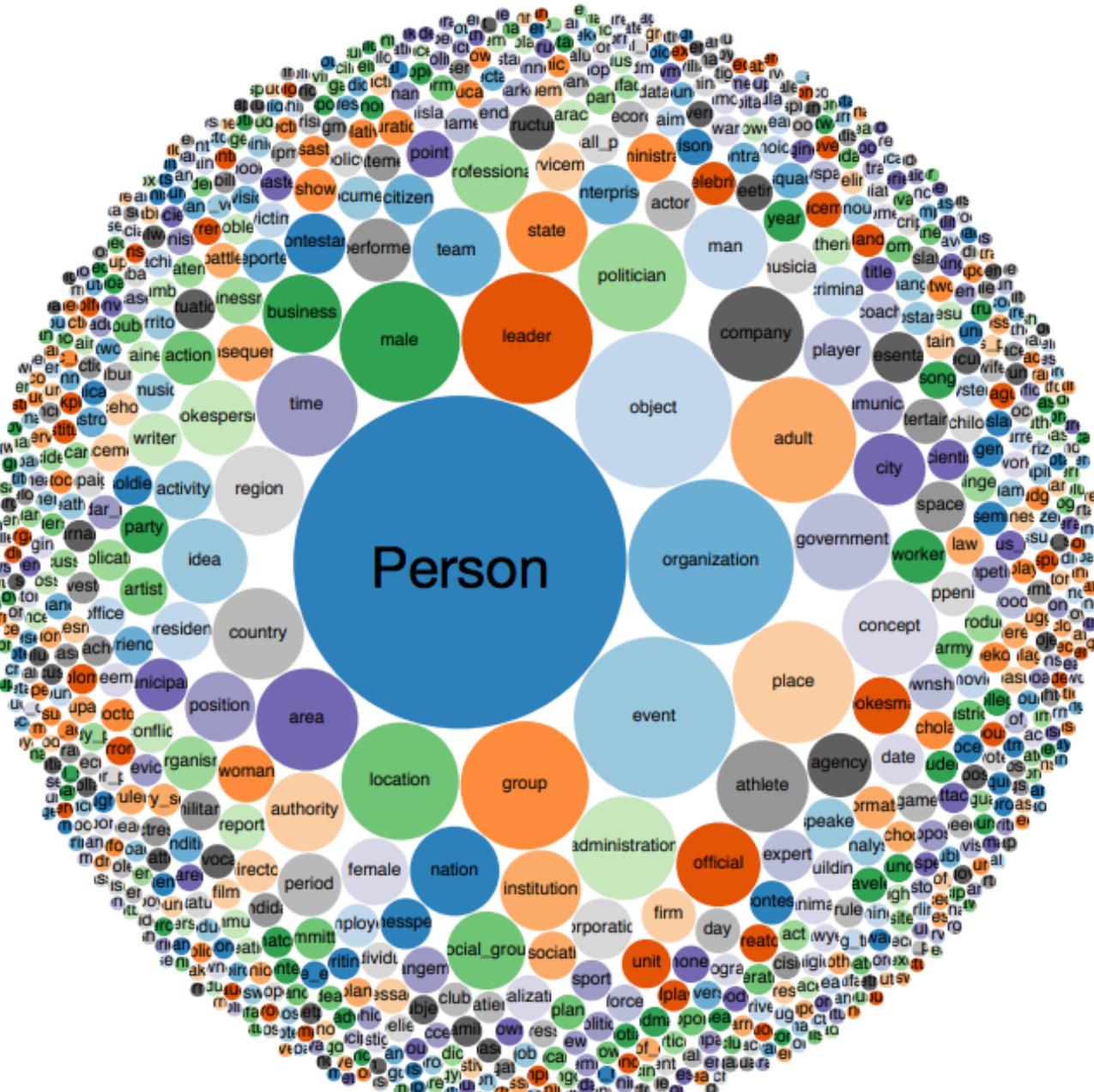
ORG

- ▶ Prediction of all Os still gets 66% accuracy on this example!
- ▶ What we really want to know: how many named entity *chunk* predictions did we get right?
- ▶ Precision: of the ones we predicted, how many are right?
- ▶ Recall: of the gold named entities, how many did we find?
- ▶ F-measure: harmonic mean of these two

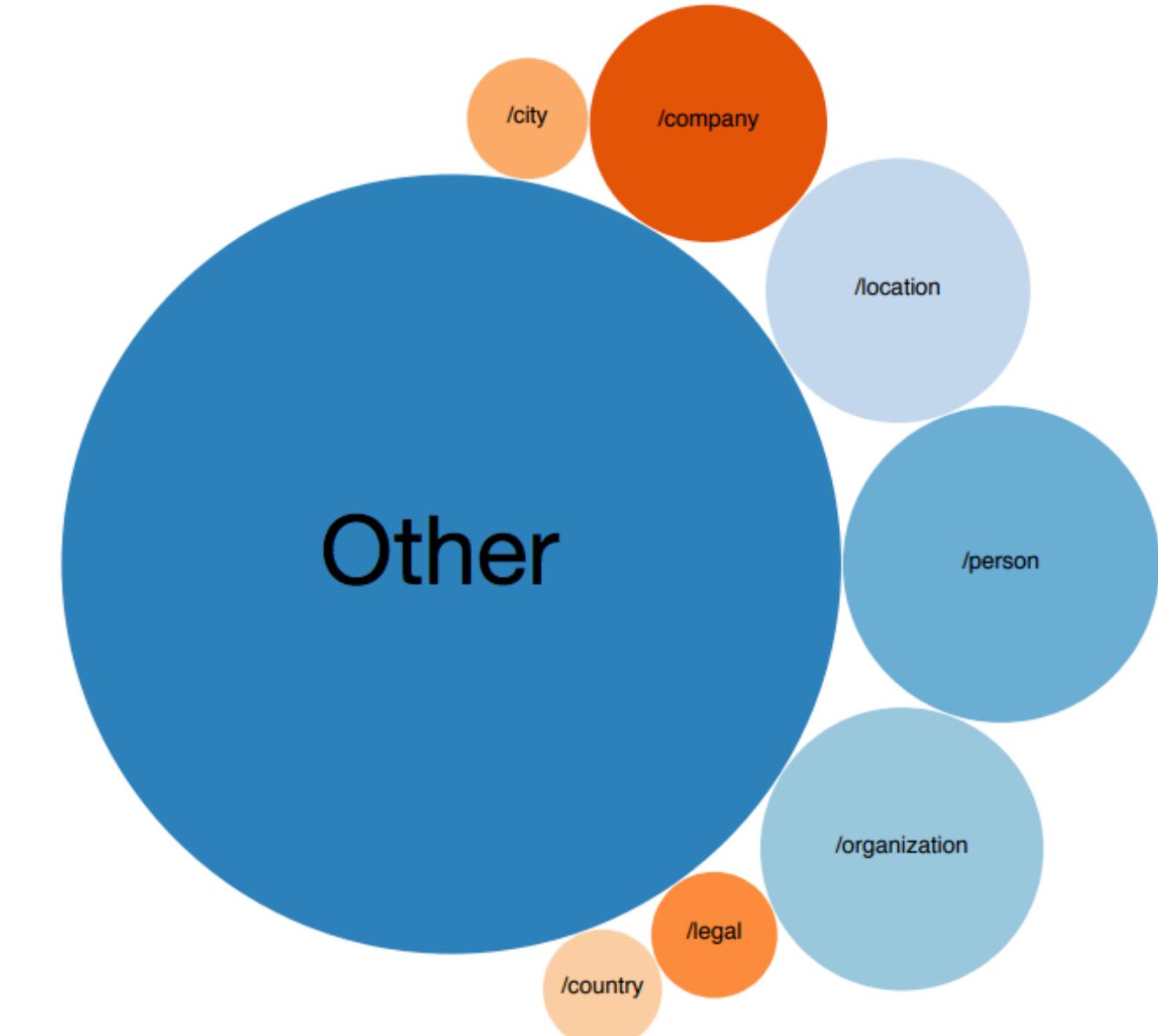
# How well do NER systems do?

	System	Resources Used	$F_1$	
+	LBJ-NER	Wikipedia, Nonlocal Features, Word-class Model	90.80	Lample et al. (2016)
-	(Suzuki and Isozaki, 2008)	Semi-supervised on 1G-word unlabeled data	89.92	LSTM-CRF (no char) <b>90.20</b>
-	(Ando and Zhang, 2005)	Semi-supervised on 27M-word unlabeled data	89.31	LSTM-CRF <b>90.94</b>
-	(Kazama and Torisawa, 2007a)	Wikipedia	88.02	S-LSTM (no char) <b>87.96</b>
-	(Krishnan and Manning, 2006)	Non-local Features	87.24	S-LSTM <b>90.33</b>
-	(Kazama and Torisawa, 2007b)	Non-local Features	87.17	BiLSTM-CRF + ELMo <b>92.2</b>
+	(Finkel et al., 2005)	Non-local Features	86.86	Peters et al. (2018)
				BERT <b>92.8</b>
				Devlin et al. (2019)

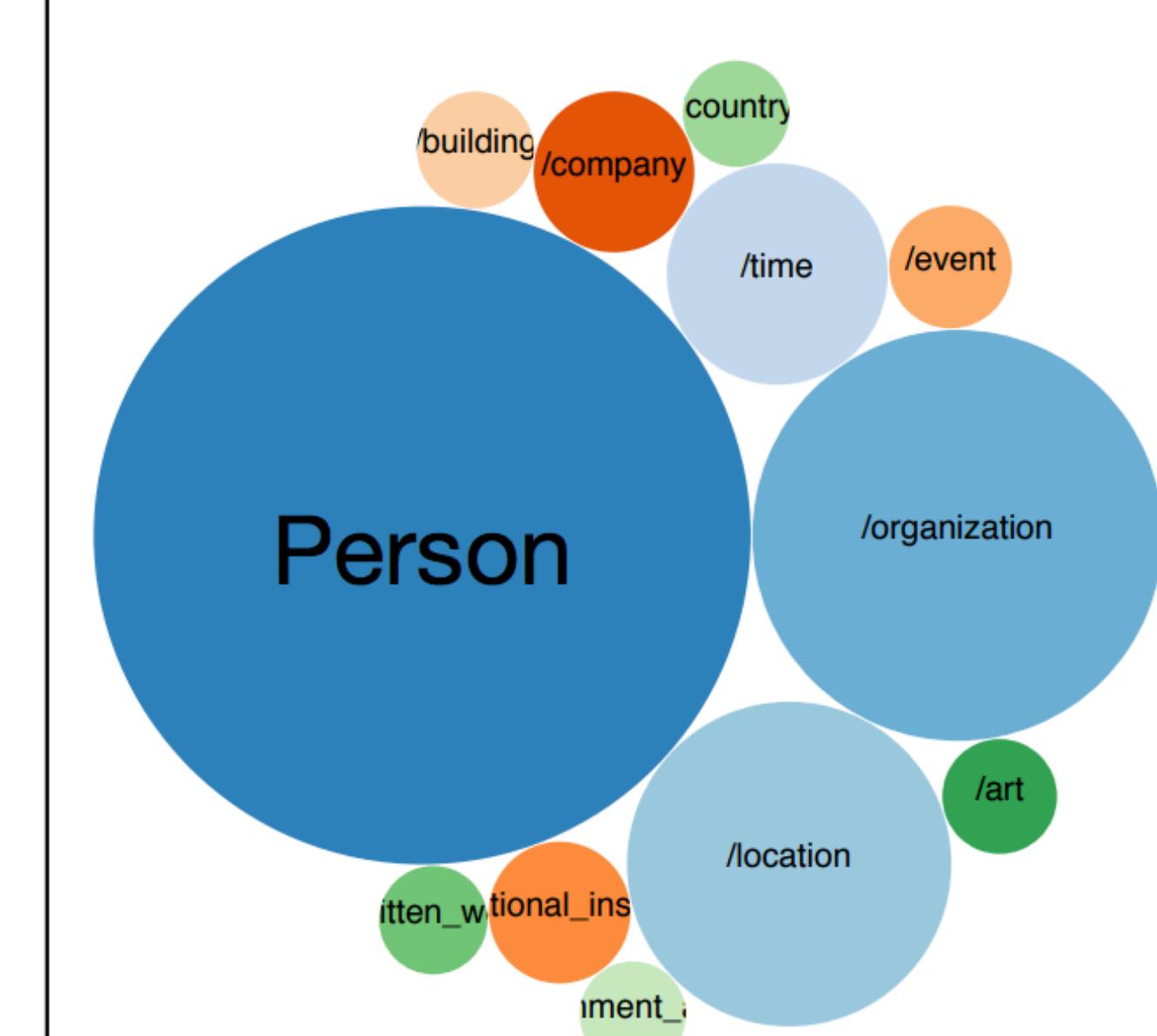
# Modern Entity Typing



### a) Our Dataset



b) OntoNotes



### c) FIGER

- More and more classes (17 → 112 → 10,000+)

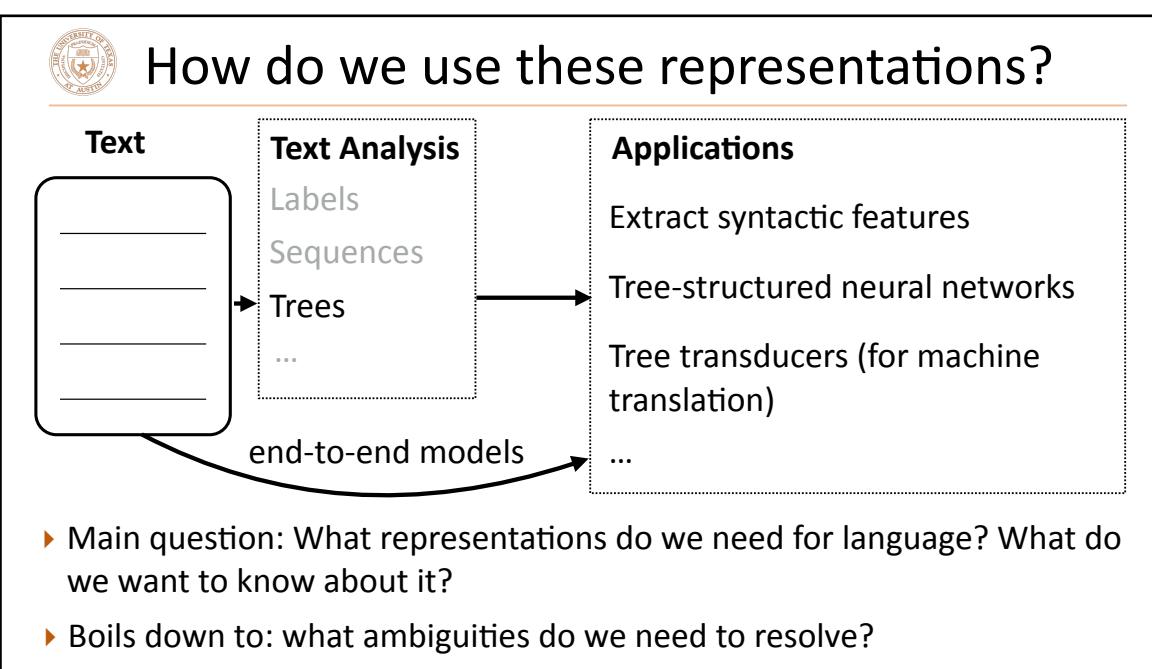
# Next Time

---

- ▶ Neural networks

# slide credits

slides that look like this



come from

CS388 given by Greg Durrett at U Texas, Austin

A screenshot of a presentation slide titled "Assignment". The slide has a dark blue header bar with the university logo. Below the header, the title "Question 2: Tagging" is displayed. To the left of the main content, there is a sidebar with several orange square bullet points:

- Six programming problems
- At least 5 must be solved
- The best two from each problem are graded.
- The sum of the 4 best scores will pass the course.
- Start early! Late assignments will not be graded.

The main content area contains the following text and equations:

- Given observations  $y_1, \dots, y_T$ , what is the most probable sequence  $x_1, \dots, x_T$  of hidden states?
- Maximum probability:
$$\max_{x_1, \dots, x_T} P(x_1, \dots, x_T | y_1, \dots, y_T)$$
- We are primarily interested in arg max:
$$\arg \max_{x_1, \dots, x_T} P(x_1, \dots, x_T | y_1, \dots, y_T) \\ = \arg \max_{x_1, \dots, x_T} \frac{P(x_1, \dots, x_T, y_1, \dots, y_T)}{P(y_1, \dots, y_T)} \\ = \arg \max_{x_1, \dots, x_T} P(x_1, \dots, x_T, y_1, \dots, y_T)$$

earlier editions of this class (ANLP), given by Tatjana Scheffler and Alexander Koller

and their use is gratefully acknowledged. I try to make any modifications obvious, but if there are errors on a slide, assume that I added them.