



Language  
Technologies  
Institute

Carnegie  
Mellon  
University

# Algorithms for NLP

CS 11-711 · Fall 2020

**Lecture 9: CRFs, neural sequence labeling**

Emma Strubell

# Announcements

- **Project 2 released today after class:** sequence labeling.
  - Due: October 16.
  - You will implement part-of-speech taggers for English and Norwegian:
    - HMM, BiLSTM, and BiLSTM-CRF.
- Friday's recitation will be an overview of P2.

# Recap

# Recap

- **HMMs**: Natural extension of Naive Bayes to sequence labeling

# Recap

- **HMMs**: Natural extension of Naive Bayes to sequence labeling
  - Hard to add rich features of the input, e.g. affixes, capitalization, ...

# Recap

- **HMMs**: Natural extension of Naive Bayes to sequence labeling
  - Hard to add rich features of the input, e.g. affixes, capitalization, ...
- Would like to train a **discriminative model**, like logistic regression, to directly model the conditional probability of labels given inputs.

# Recap

- **HMMs**: Natural extension of Naive Bayes to sequence labeling
  - Hard to add rich features of the input, e.g. affixes, capitalization, ...
- Would like to train a **discriminative model**, like logistic regression, to directly model the conditional probability of labels given inputs.
- Logistic regression (MEMMs) suffer from the **label bias problem**.

# Recap

- **HMMs**: Natural extension of Naive Bayes to sequence labeling
  - Hard to add rich features of the input, e.g. affixes, capitalization, ...
- Would like to train a **discriminative model**, like logistic regression, to directly model the conditional probability of labels given inputs.
  - Logistic regression (MEMMs) suffer from the **label bias problem**.
- Solution: **linear-chain CRFs**.

# Conditional random fields (CRFs)

- **Linear-chain CRFs**: Globally-normalized discriminative sequence labeling models!

# Conditional random fields (CRFs)

- **Linear-chain CRFs**: Globally-normalized discriminative sequence labeling models!

$$P(\mathbf{y} \mid \mathbf{w}) = \frac{\exp(\Psi(\mathbf{w}, \mathbf{y}))}{\sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{w})} \exp(\Psi(\mathbf{w}, \mathbf{y}'))}$$

# Conditional random fields (CRFs)

- **Linear-chain CRFs**: Globally-normalized discriminative sequence labeling models!

$$P(\mathbf{y} \mid \mathbf{w}) = \frac{\exp(\Psi(\mathbf{w}, \mathbf{y}))}{\sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{w})} \exp(\Psi(\mathbf{w}, \mathbf{y}'))}$$

$\mathbf{y}'$  is an entire sequence

# Conditional random fields (CRFs)

- **Linear-chain CRFs**: Globally-normalized discriminative sequence labeling models!

$$P(\mathbf{y} \mid \mathbf{w}) = \frac{\exp(\Psi(\mathbf{w}, \mathbf{y}))}{\sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{w})} \exp(\Psi(\mathbf{w}, \mathbf{y}'))}$$

$\mathbf{y}'$  is an entire sequence

$$\Psi(\mathbf{w}, \mathbf{y}) = \sum_{m=1}^{M+1} \psi(\mathbf{w}, y_m, y_{m-1}, m)$$

# Conditional random fields (CRFs)

- **Linear-chain CRFs**: Globally-normalized discriminative sequence labeling models!

$$P(\mathbf{y} \mid \mathbf{w}) = \frac{\exp(\Psi(\mathbf{w}, \mathbf{y}))}{\sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{w})} \exp(\Psi(\mathbf{w}, \mathbf{y}'))}$$

$\mathbf{y}'$  is an entire sequence

$$\Psi(\mathbf{w}, \mathbf{y}) = \sum_{m=1}^{M+1} \psi(\mathbf{w}, y_m, y_{m-1}, m)$$

decompose into local scores

# Conditional random fields (CRFs)

- **Linear-chain CRFs**: Globally-normalized discriminative sequence labeling models!

$$P(\mathbf{y} \mid \mathbf{w}) = \frac{\exp(\Psi(\mathbf{w}, \mathbf{y}))}{\sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{w})} \exp(\Psi(\mathbf{w}, \mathbf{y}'))}$$

$\mathbf{y}'$  is an entire sequence

$$\Psi(\mathbf{w}, \mathbf{y}) = \sum_{m=1}^{M+1} \psi(\mathbf{w}, y_m, y_{m-1}, m)$$

decompose into local scores

$$= \sum_{m=1}^{M+1} \theta \cdot \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m)$$

# Conditional random fields (CRFs)

- **Linear-chain CRFs**: Globally-normalized discriminative sequence labeling models!

$$\Psi(\mathbf{w}, \mathbf{y}) = \sum_{m=1}^{M+1} \theta \cdot \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m)$$

# Conditional random fields (CRFs)

- **Linear-chain CRFs**: Globally-normalized discriminative sequence labeling models!

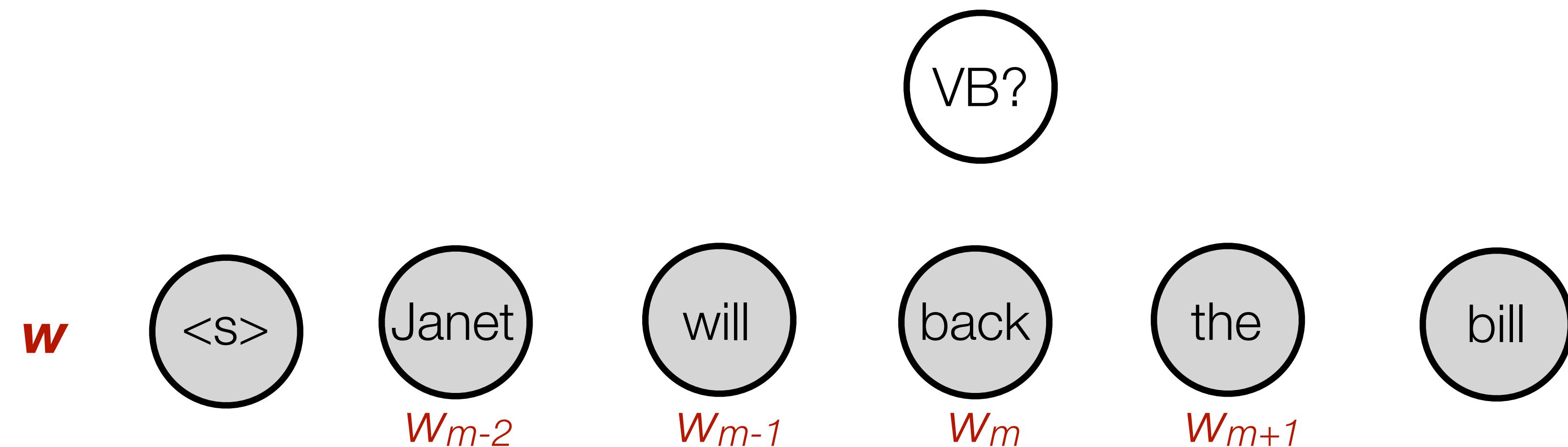
$$\Psi(\mathbf{w}, \mathbf{y}) = \sum_{m=1}^{M+1} \theta \cdot \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m)$$



# Conditional random fields (CRFs)

- **Linear-chain CRFs**: Globally-normalized discriminative sequence labeling models!

$$\Psi(\mathbf{w}, \mathbf{y}) = \sum_{m=1}^{M+1} \theta \cdot \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m)$$

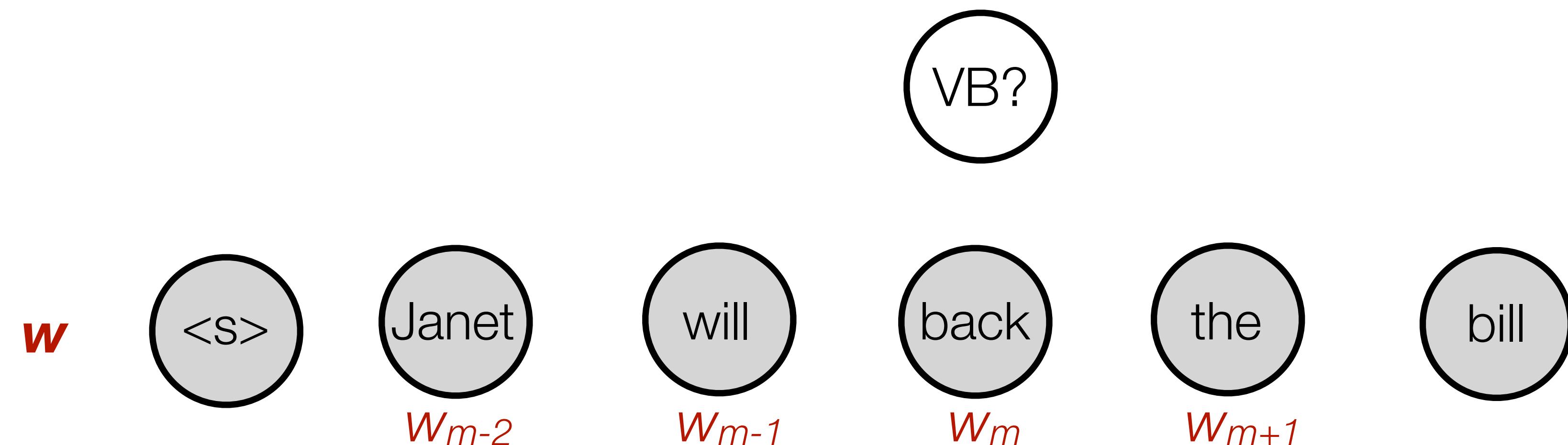


# Conditional random fields (CRFs)

- **Linear-chain CRFs**: Globally-normalized discriminative sequence labeling models!

$$\Psi(\mathbf{w}, \mathbf{y}) = \sum_{m=1}^{M+1} \theta \cdot \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m)$$

1	0	...	1	...	0	...	1	0	0	0	0	0	...	0
---	---	-----	---	-----	---	-----	---	---	---	---	---	---	-----	---



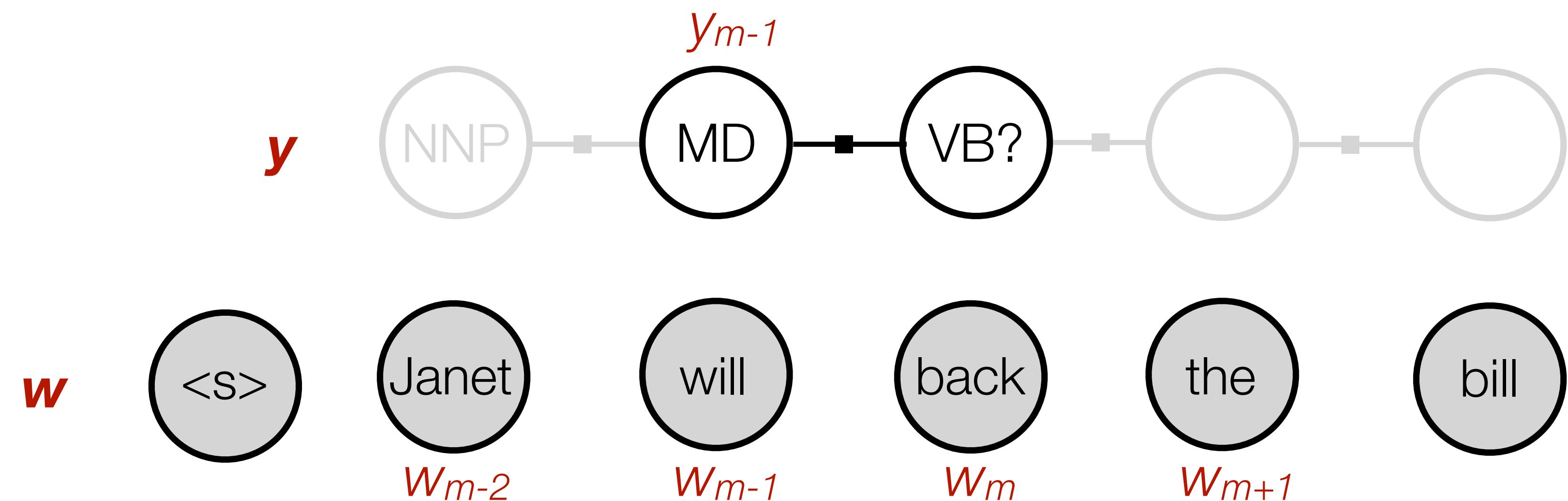
# Conditional random fields (CRFs)

- **Linear-chain CRFs**: Globally-normalized discriminative sequence labeling models!

$$\Psi(\mathbf{w}, \mathbf{y}) = \sum_{m=1}^{M+1} \theta \cdot \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m)$$

1	0	...	1	...	0	...	1	0	0	0	0	0	...	0
---	---	-----	---	-----	---	-----	---	---	---	---	---	---	-----	---

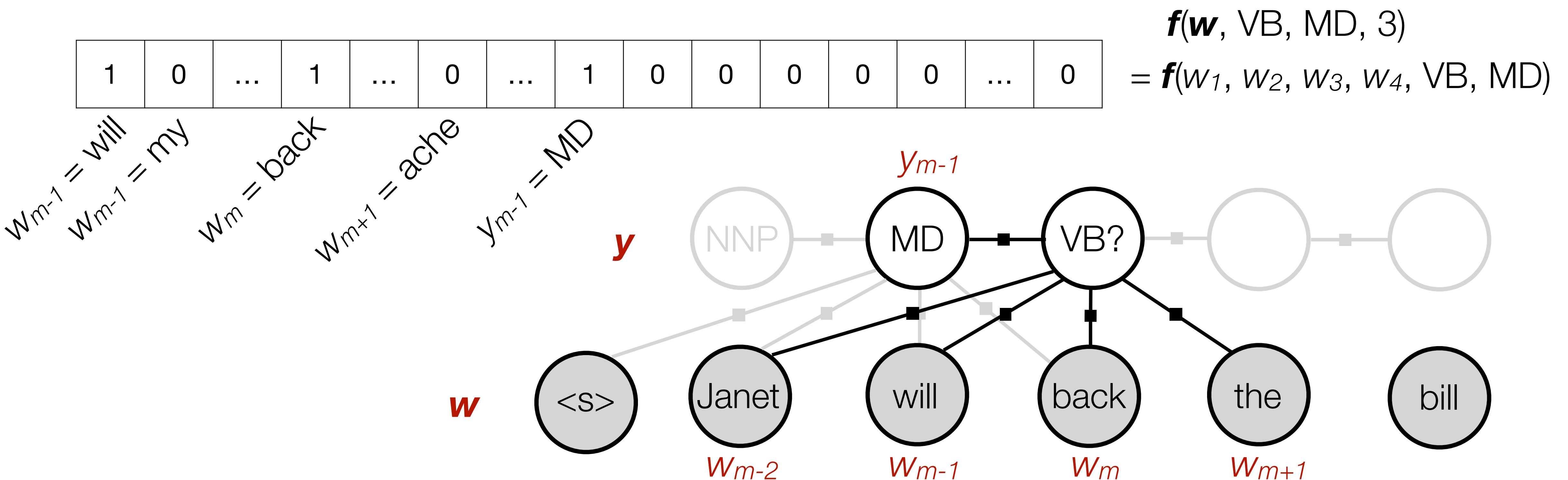
$\mathbf{f}(\mathbf{w}, \text{VB}, \text{MD}, 3)$



# Conditional random fields (CRFs)

- **Linear-chain CRFs**: Globally-normalized discriminative sequence labeling models!

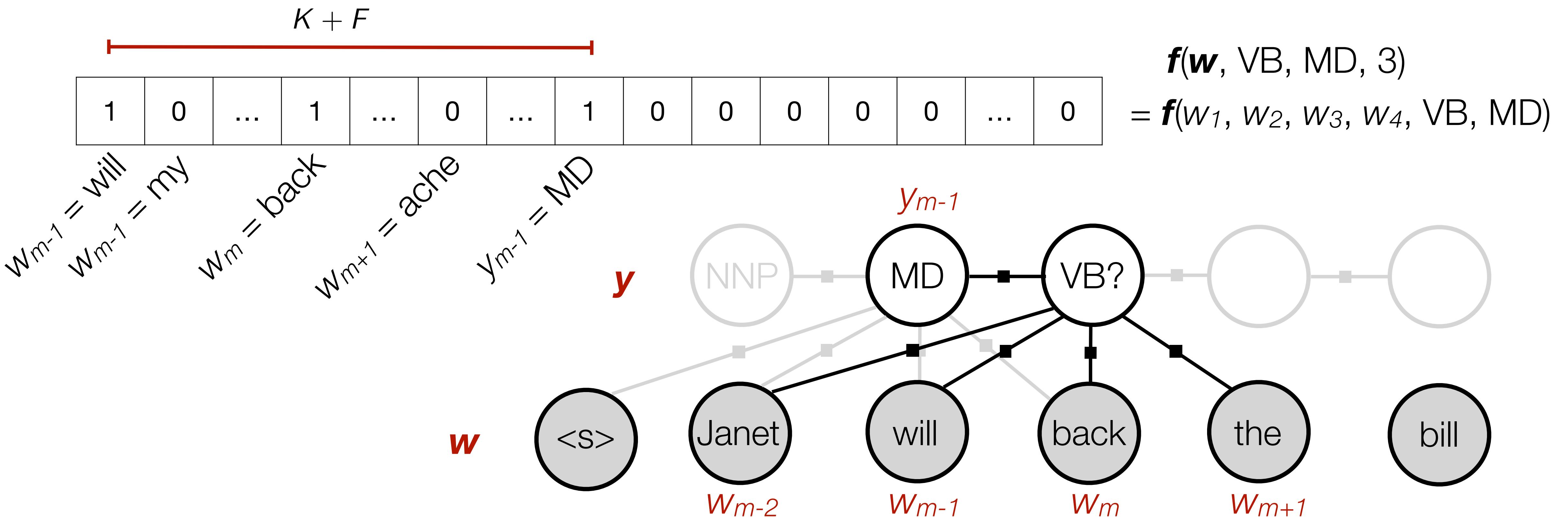
$$\Psi(\mathbf{w}, \mathbf{y}) = \sum_{m=1}^{M+1} \theta \cdot \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m)$$



# Conditional random fields (CRFs)

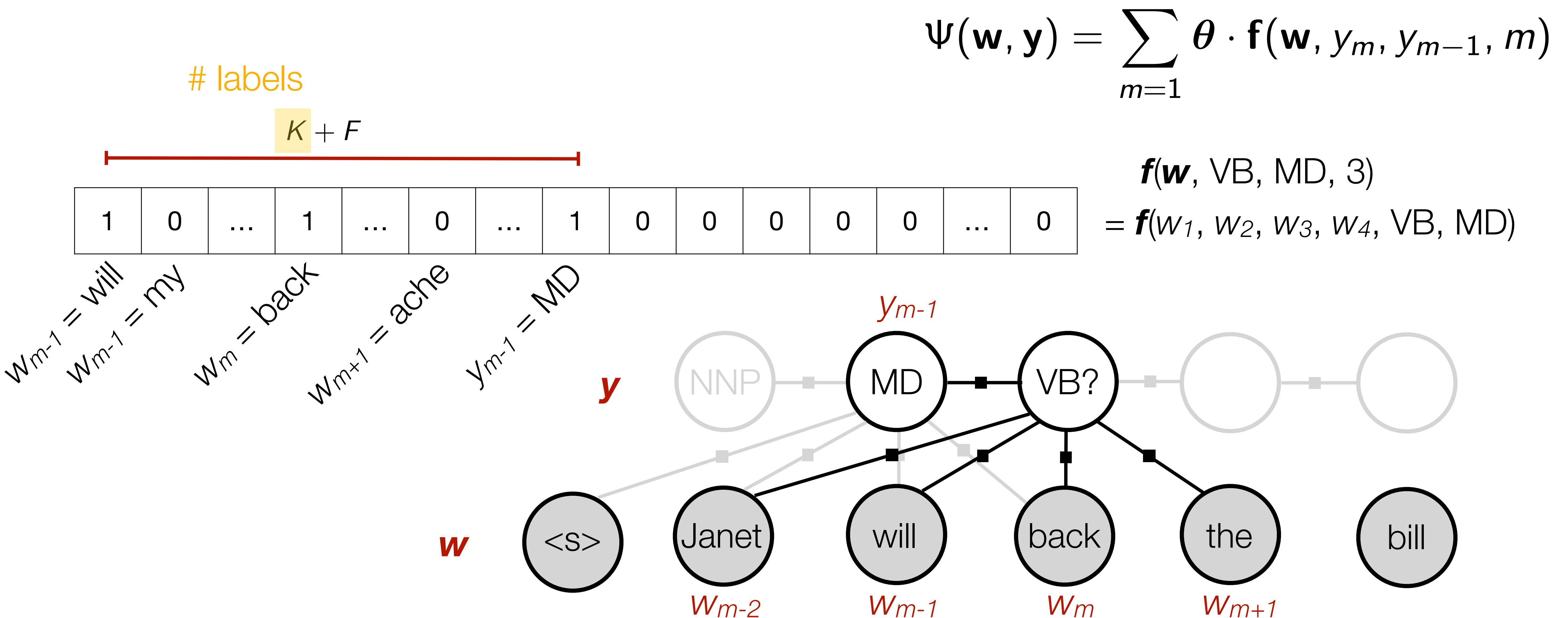
- **Linear-chain CRFs**: Globally-normalized discriminative sequence labeling models!

$$\Psi(\mathbf{w}, \mathbf{y}) = \sum_{m=1}^{M+1} \theta \cdot \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m)$$



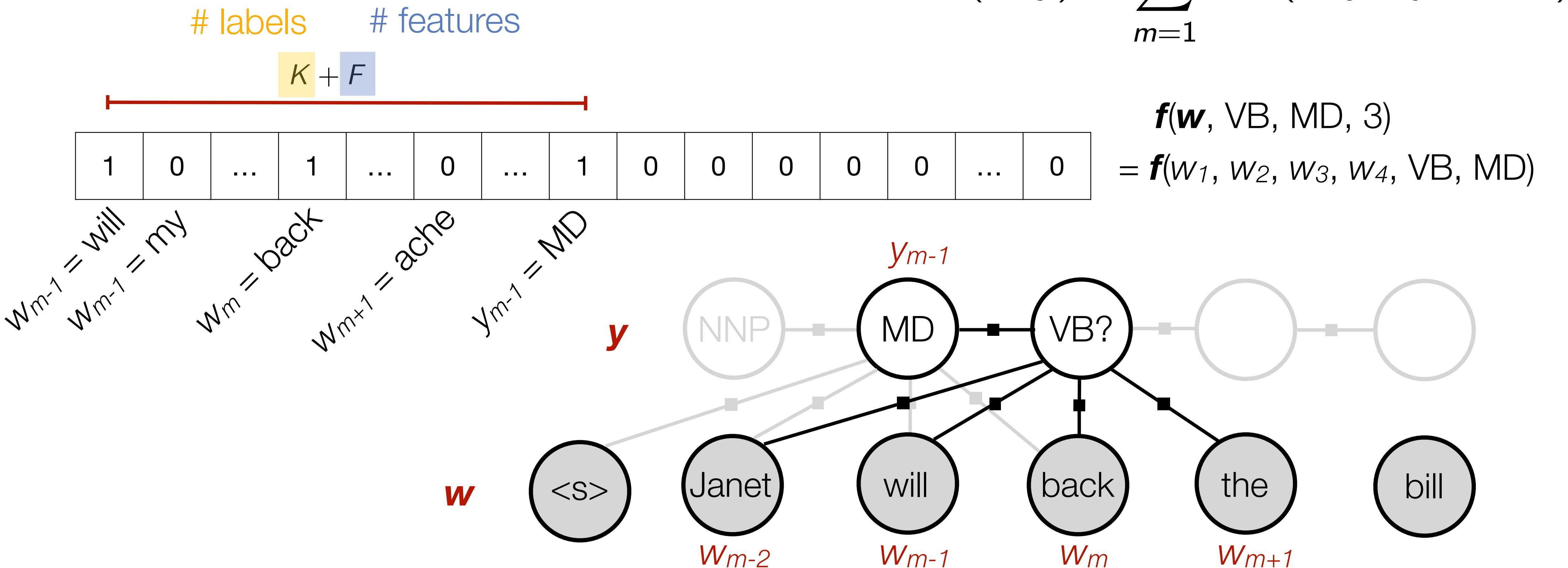
# Conditional random fields (CRFs)

- **Linear-chain CRFs**: Globally-normalized discriminative sequence labeling models!



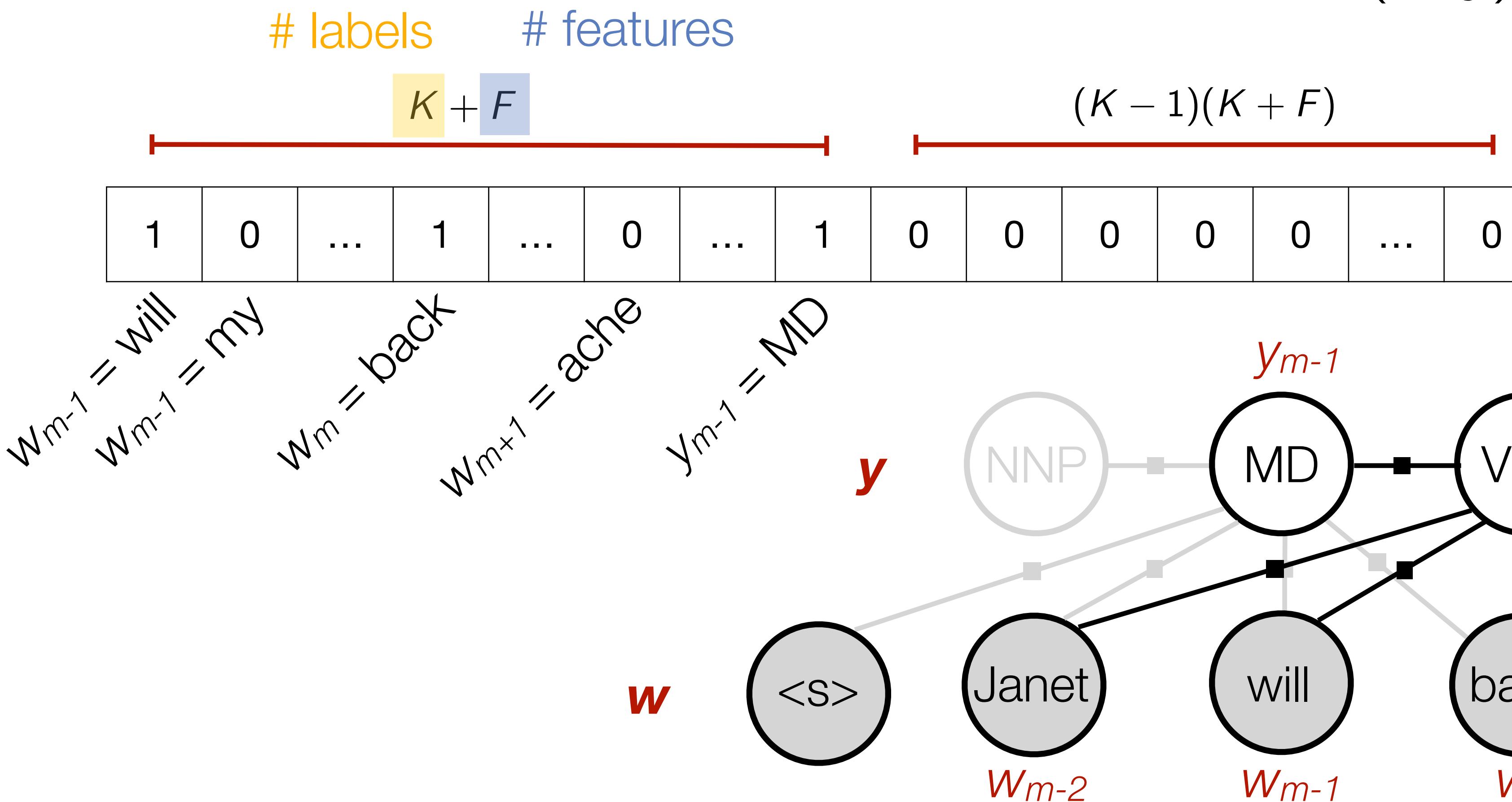
# Conditional random fields (CRFs)

- **Linear-chain CRFs**: Globally-normalized discriminative sequence labeling models!



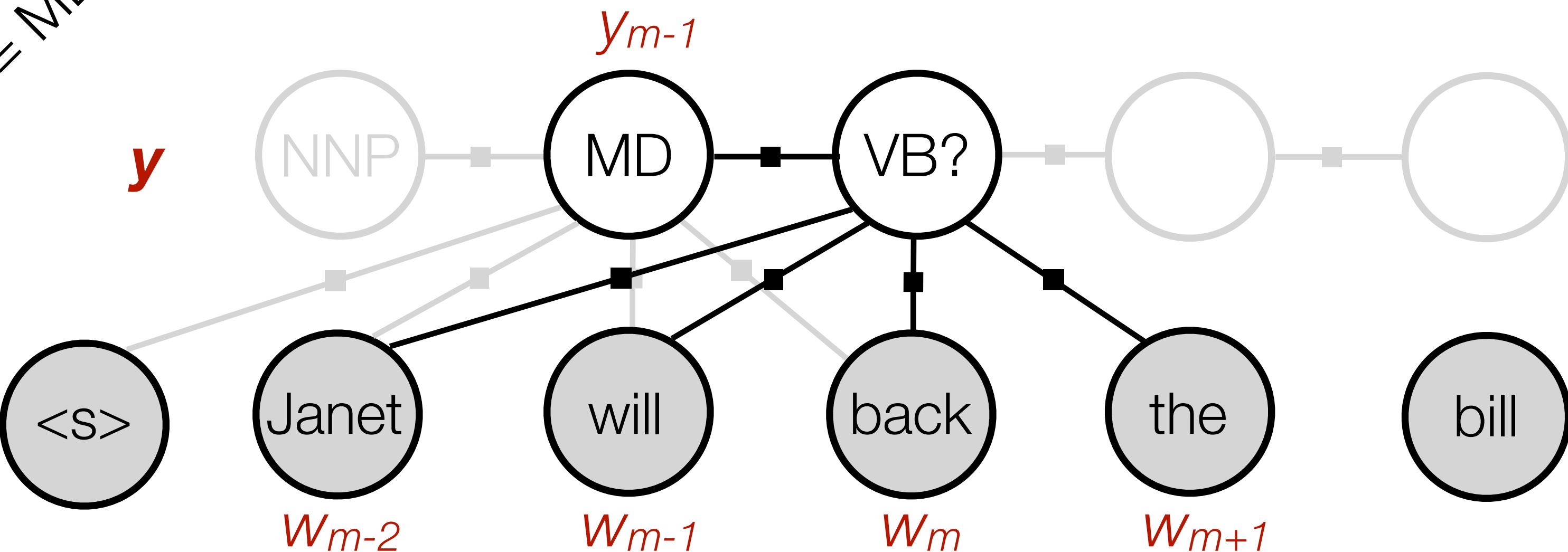
# Conditional random fields (CRFs)

- **Linear-chain CRFs**: Globally-normalized discriminative sequence labeling models!



$$\Psi(\mathbf{w}, \mathbf{y}) = \sum_{m=1}^{M+1} \theta \cdot \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m)$$

$$\begin{aligned} \mathbf{f}(\mathbf{w}, \text{VB}, \text{MD}, 3) \\ = \mathbf{f}(w_1, w_2, w_3, w_4, \text{VB}, \text{MD}) \end{aligned}$$



# Conditional random fields (CRFs)

- **Decoding:**

$$P(\mathbf{y} \mid \mathbf{w}) = \frac{\exp(\Psi(\mathbf{w}, \mathbf{y}))}{\sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{w})} \exp(\Psi(\mathbf{w}, \mathbf{y}'))}$$

# Conditional random fields (CRFs)

- **Decoding:** Direct application of Viterbi!

$$P(\mathbf{y} \mid \mathbf{w}) = \frac{\exp(\Psi(\mathbf{w}, \mathbf{y}))}{\sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{w})} \exp(\Psi(\mathbf{w}, \mathbf{y}'))}$$

# Conditional random fields (CRFs)

- **Decoding:** Direct application of Viterbi!

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} \log P(\mathbf{y} \mid \mathbf{w})$$

$$P(\mathbf{y} \mid \mathbf{w}) = \frac{\exp(\Psi(\mathbf{w}, \mathbf{y}))}{\sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{w})} \exp(\Psi(\mathbf{w}, \mathbf{y}'))}$$

# Conditional random fields (CRFs)

- **Decoding:** Direct application of Viterbi!

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} \log P(\mathbf{y} \mid \mathbf{w})$$

$$= \operatorname{argmax}_{\mathbf{y}} \Psi(\mathbf{y}, \mathbf{w}) - \log \sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{w})} \exp \Psi(\mathbf{y}', \mathbf{w})$$

$$P(\mathbf{y} \mid \mathbf{w}) = \frac{\exp(\Psi(\mathbf{w}, \mathbf{y}))}{\sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{w})} \exp(\Psi(\mathbf{w}, \mathbf{y}'))}$$



# Conditional random fields (CRFs)

- **Decoding:** Direct application of Viterbi!

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} \log P(\mathbf{y} | \mathbf{w})$$

$$= \operatorname{argmax}_{\mathbf{y}} \Psi(\mathbf{y}, \mathbf{w}) - \log \sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{w})} \exp \Psi(\mathbf{y}', \mathbf{w})$$

$$P(\mathbf{y} | \mathbf{w}) = \frac{\exp(\Psi(\mathbf{w}, \mathbf{y}))}{\sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{w})} \exp(\Psi(\mathbf{w}, \mathbf{y}'))}$$

# Conditional random fields (CRFs)

- **Decoding:** Direct application of Viterbi!

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} \log P(\mathbf{y} | \mathbf{w})$$

$$= \operatorname{argmax}_{\mathbf{y}} \Psi(\mathbf{y}, \mathbf{w}) - \log \sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{w})} \exp \Psi(\mathbf{y}', \mathbf{w})$$

$$= \operatorname{argmax}_{\mathbf{y}} \Psi(\mathbf{y}, \mathbf{w}) \quad \text{same for all settings of } \mathbf{y}$$

$$P(\mathbf{y} | \mathbf{w}) = \frac{\exp(\Psi(\mathbf{w}, \mathbf{y}))}{\sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{w})} \exp(\Psi(\mathbf{w}, \mathbf{y}'))}$$

# Conditional random fields (CRFs)

- **Decoding:** Direct application of Viterbi!

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} \log P(\mathbf{y} | \mathbf{w})$$

$$= \operatorname{argmax}_{\mathbf{y}} \Psi(\mathbf{y}, \mathbf{w}) - \log \sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{w})} \exp \Psi(\mathbf{y}', \mathbf{w})$$

$$= \operatorname{argmax}_{\mathbf{y}} \Psi(\mathbf{y}, \mathbf{w}) \quad \text{same for all settings of } \mathbf{y}$$

$$= \operatorname{argmax}_{\mathbf{y}} \sum_{m=1}^{M+1} \psi(\mathbf{w}, y_m, y_{m-1}, m)$$

$$P(\mathbf{y} | \mathbf{w}) = \frac{\exp(\Psi(\mathbf{w}, \mathbf{y}))}{\sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{w})} \exp(\Psi(\mathbf{w}, \mathbf{y}'))}$$

# Conditional random fields (CRFs)

- **Decoding:** Direct application of Viterbi!

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} \log P(\mathbf{y} | \mathbf{w})$$

$$= \operatorname{argmax}_{\mathbf{y}} \Psi(\mathbf{y}, \mathbf{w}) - \log \sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{w})} \exp \Psi(\mathbf{y}', \mathbf{w})$$

$$= \operatorname{argmax}_{\mathbf{y}} \Psi(\mathbf{y}, \mathbf{w}) \quad \text{same for all settings of } \mathbf{y}$$

$$= \operatorname{argmax}_{\mathbf{y}} \sum_{m=1}^{M+1} \boxed{\psi(\mathbf{w}, y_m, y_{m-1}, m) = s_m(y_m, y_{m-1})}$$

$$P(\mathbf{y} | \mathbf{w}) = \frac{\exp(\Psi(\mathbf{w}, \mathbf{y}))}{\sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{w})} \exp(\Psi(\mathbf{w}, \mathbf{y}'))}$$

# Conditional random fields (CRFs)

- **Decoding:** Direct application of Viterbi!

$$\hat{\mathbf{y}} = \underset{\mathbf{y}}{\operatorname{argmax}} \sum_{m=1}^{M+1} \psi(\mathbf{w}, y_m, y_{m-1}, m) = s_m(y_m, y_{m-1})$$

# Conditional random fields (CRFs)

- **Decoding:** Direct application of Viterbi!

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} \sum_{m=1}^{M+1} \psi(\mathbf{w}, y_m, y_{m-1}, m) = s_m(y_m, y_{m-1})$$

$$\max_{\mathbf{y}_{1:M}} \Psi(\mathbf{w}, \mathbf{y}_{1:M}) = \max_{\mathbf{y}_{1:M}} \sum_{m=1}^{M+1} s_m(y_m, y_{m-1})$$

# Conditional random fields (CRFs)

- **Decoding:** Direct application of Viterbi!

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} \sum_{m=1}^{M+1} \psi(\mathbf{w}, y_m, y_{m-1}, m) = s_m(y_m, y_{m-1})$$

score of best tag sequence  
of length M

$$\max_{\mathbf{y}_{1:M}} \Psi(\mathbf{w}, \mathbf{y}_{1:M}) = \max_{\mathbf{y}_{1:M}} \sum_{m=1}^{M+1} s_m(y_m, y_{m-1})$$

# Conditional random fields (CRFs)

- **Decoding:** Direct application of Viterbi!

$$\hat{\mathbf{y}} = \underset{\mathbf{y}}{\operatorname{argmax}} \sum_{m=1}^{M+1} \psi(\mathbf{w}, y_m, y_{m-1}, m) = s_m(y_m, y_{m-1})$$

score of best tag sequence  
of length M

$$\begin{aligned} \max_{\mathbf{y}_{1:M}} \Psi(\mathbf{w}, \mathbf{y}_{1:M}) &= \max_{\mathbf{y}_{1:M}} \sum_{m=1}^{M+1} s_m(y_m, y_{m-1}) \\ &= \max_{\mathbf{y}_M} s_{M+1}(\langle /s \rangle, y_M) + \max_{\mathbf{y}_{1:M-1}} \sum_{m=1}^M s_m(y_m, y_{m-1}) \end{aligned}$$

# Conditional random fields (CRFs)

- **Decoding:** Direct application of Viterbi!

$$\hat{\mathbf{y}} = \underset{\mathbf{y}}{\operatorname{argmax}} \sum_{m=1}^{M+1} \psi(\mathbf{w}, y_m, y_{m-1}, m) = s_m(y_m, y_{m-1})$$

score of best tag sequence  
of length M

$$\begin{aligned} \max_{\mathbf{y}_{1:M}} \Psi(\mathbf{w}, \mathbf{y}_{1:M}) &= \max_{\mathbf{y}_{1:M}} \sum_{m=1}^{M+1} s_m(y_m, y_{m-1}) \\ &= \max_{\mathbf{y}_M} s_{M+1}(\langle /s \rangle, y_M) + \max_{\mathbf{y}_{1:M-1}} \sum_{m=1}^M s_m(y_m, y_{m-1}) \end{aligned}$$

# Conditional random fields (CRFs)

- **Decoding:** Direct application of Viterbi!

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} \sum_{m=1}^{M+1} \psi(\mathbf{w}, y_m, y_{m-1}, m) = s_m(y_m, y_{m-1})$$

score of best tag sequence  
of length M

$$\begin{aligned} \max_{\mathbf{y}_{1:M}} \Psi(\mathbf{w}, \mathbf{y}_{1:M}) &= \max_{\mathbf{y}_{1:M}} \sum_{m=1}^{M+1} s_m(y_m, y_{m-1}) \\ &= \max_{\mathbf{y}_M} s_{M+1}(\langle /s \rangle, y_M) + \max_{\mathbf{y}_{1:M-1}} \sum_{m=1}^M s_m(y_m, y_{m-1}) \end{aligned}$$

score of  
most probable extension  $y_M$

# Conditional random fields (CRFs)

- **Decoding:** Direct application of Viterbi!

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} \sum_{m=1}^{M+1} \psi(\mathbf{w}, y_m, y_{m-1}, m) = s_m(y_m, y_{m-1})$$

score of best tag sequence  
of length M

$$\begin{aligned} \max_{\mathbf{y}_{1:M}} \Psi(\mathbf{w}, \mathbf{y}_{1:M}) &= \max_{\mathbf{y}_{1:M}} \sum_{m=1}^{M+1} s_m(y_m, y_{m-1}) \\ &= \max_{\mathbf{y}_M} s_{M+1}(\langle /s \rangle, y_M) + \max_{\mathbf{y}_{1:M-1}} \sum_{m=1}^M s_m(y_m, y_{m-1}) \end{aligned}$$

same subproblem!

score of most probable extension  $y_M$

# Conditional random fields (CRFs)

- **Decoding:** Direct application of Viterbi!

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} \sum_{m=1}^{M+1} \psi(\mathbf{w}, y_m, y_{m-1}, m) = s_m(y_m, y_{m-1})$$

score of best tag sequence  
of length M

$$\begin{aligned} \max_{\mathbf{y}_{1:M}} \Psi(\mathbf{w}, \mathbf{y}_{1:M}) &= \max_{\mathbf{y}_{1:M}} \sum_{m=1}^{M+1} s_m(y_m, y_{m-1}) \\ &= \max_{\mathbf{y}_M} s_{M+1}(\langle /s \rangle, y_M) + \max_{\mathbf{y}_{1:M-1}} \sum_{m=1}^M s_m(y_m, y_{m-1}) \end{aligned}$$

same subproblem!

score of most probable extension  $y_M$

# Learning in CRFs

# Learning in CRFs

- As with logistic regression, weights  $\theta$  are learned by minimizing negative log likelihood:

$$\ell = - \sum_{i=1}^N \log P(\mathbf{y}^{(i)} \mid \mathbf{w}^{(i)}; \theta)$$

# Learning in CRFs

- As with logistic regression, weights  $\theta$  are learned by minimizing negative log likelihood:

$$\begin{aligned}\ell &= - \sum_{i=1}^N \log P(\mathbf{y}^{(i)} \mid \mathbf{w}^{(i)}; \theta) \\ &= - \sum_{i=1}^N \theta \cdot \mathbf{f}(\mathbf{w}^{(i)}, \mathbf{y}^{(i)}) + \log \sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{w}^{(i)})} \exp(\theta \cdot \mathbf{f}(\mathbf{w}^{(i)}, \mathbf{y}'))\end{aligned}$$

# Learning in CRFs

- As with logistic regression, weights  $\theta$  are learned by minimizing negative log likelihood:

$$\begin{aligned}\ell &= - \sum_{i=1}^N \log P(\mathbf{y}^{(i)} \mid \mathbf{w}^{(i)}; \theta) \\ &= - \sum_{i=1}^N \theta \cdot \mathbf{f}(\mathbf{w}^{(i)}, \mathbf{y}^{(i)}) + \log \sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{w}^{(i)})} \exp(\theta \cdot \mathbf{f}(\mathbf{w}^{(i)}, \mathbf{y}'))\end{aligned}$$

sum over all possible labelings

# Learning in CRFs

- As with logistic regression, weights  $\theta$  are learned by minimizing negative log likelihood:

$$\begin{aligned}\ell &= - \sum_{i=1}^N \log P(\mathbf{y}^{(i)} \mid \mathbf{w}^{(i)}; \theta) \\ &= - \sum_{i=1}^N \theta \cdot \mathbf{f}(\mathbf{w}^{(i)}, \mathbf{y}^{(i)}) + \log \sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{w}^{(i)})} \exp(\theta \cdot \mathbf{f}(\mathbf{w}^{(i)}, \mathbf{y}'))\end{aligned}$$

sum over all possible labelings

- Can be computed efficiently using **forward algorithm**.

# Learning in CRFs

# Learning in CRFs

- Likelihood can be computed efficiently using **forward algorithm**.

# Learning in CRFs

- Likelihood can be computed efficiently using **forward algorithm**. Define:

$$\alpha_m(y_m) = \sum_{\mathbf{y}_{1:m-1}} \exp \sum_{n=1}^m s_n(y_n, y_{n-1})$$

# Learning in CRFs

- Likelihood can be computed efficiently using **forward algorithm**. Define:

$$\alpha_m(y_m) = \sum_{\mathbf{y}_{1:m-1}} \exp \sum_{n=1}^m s_n(y_n, y_{n-1}) = \sum_{\mathbf{y}_{1:m-1}} \prod_{n=1}^m \exp s_n(y_n, y_{n-1})$$

# Learning in CRFs

- Likelihood can be computed efficiently using **forward algorithm**. Define:

$$\alpha_m(y_m) = \sum_{\mathbf{y}_{1:m-1}} \exp \sum_{n=1}^m s_n(y_n, y_{n-1}) = \sum_{\mathbf{y}_{1:m-1}} \prod_{n=1}^m \exp s_n(y_n, y_{n-1})$$

- As in decoding / Viterbi, can be decomposed into recursive substructure:

$$= \sum_{y_{m-1}} (\exp s_m(y_m, y_{m-1})) \sum_{\mathbf{y}_{1:m-2}} \prod_{n=1}^{m-1} \exp s_n(y_n, y_{n-1})$$

# Learning in CRFs

- Likelihood can be computed efficiently using **forward algorithm**. Define:

$$\alpha_m(y_m) = \sum_{\mathbf{y}_{1:m-1}} \exp \sum_{n=1}^m s_n(y_n, y_{n-1}) = \sum_{\mathbf{y}_{1:m-1}} \prod_{n=1}^m \exp s_n(y_n, y_{n-1})$$

- As in decoding / Viterbi, can be decomposed into recursive substructure:

$$\begin{aligned} &= \sum_{y_{m-1}} (\exp s_m(y_m, y_{m-1})) \sum_{\mathbf{y}_{1:m-2}} \prod_{n=1}^{m-1} \exp s_n(y_n, y_{n-1}) \\ &= \sum_{y_{m-1}} (\exp s_m(y_m, y_{m-1})) \times \alpha_{m-1}(y_{m-1}) \end{aligned}$$

# Learning in CRFs

- Likelihood can be computed efficiently using **forward algorithm**. Define:

$$\alpha_m(y_m) = \sum_{\mathbf{y}_{1:m-1}} \exp \sum_{n=1}^m s_n(y_n, y_{n-1}) = \sum_{\mathbf{y}_{1:m-1}} \prod_{n=1}^m \exp s_n(y_n, y_{n-1})$$

- As in decoding / Viterbi, can be decomposed into recursive substructure:

$$= \sum_{y_{m-1}} (\exp s_m(y_m, y_{m-1})) \sum_{\mathbf{y}_{1:m-2}} \prod_{n=1}^{m-1} \exp s_n(y_n, y_{n-1})$$

$$= \sum_{y_{m-1}} (\exp s_m(y_m, y_{m-1})) \times \alpha_{m-1}(y_{m-1})$$

sum instead of max

# Learning in CRFs

- Likelihood can be computed efficiently using **forward algorithm**. Define:

$$\alpha_m(y_m) = \sum_{\mathbf{y}_{1:m-1}} \exp \sum_{n=1}^m s_n(y_n, y_{n-1}) = \sum_{\mathbf{y}_{1:m-1}} \prod_{n=1}^m \exp s_n(y_n, y_{n-1})$$

- As in decoding / Viterbi, can be decomposed into recursive substructure:

$$= \sum_{y_{m-1}} (\exp s_m(y_m, y_{m-1})) \sum_{\mathbf{y}_{1:m-2}} \prod_{n=1}^{m-1} \exp s_n(y_n, y_{n-1})$$

$$= \sum_{y_{m-1}} (\exp s_m(y_m, y_{m-1})) \times \alpha_{m-1}(y_{m-1})$$

sum instead of max

Viterbi is a special case of the *max-product algorithm*,  
forward is a special case of the *sum-product algorithm*.

# Learning in CRFs

- Likelihood can be computed efficiently using **forward algorithm**. Define:

$$\alpha_m(y_m) = \sum_{\mathbf{y}_{1:m-1}} \exp \sum_{n=1}^m s_n(y_n, y_{n-1}) = \sum_{\mathbf{y}_{1:m-1}} \prod_{n=1}^m \exp s_n(y_n, y_{n-1})$$

- As in decoding / Viterbi, can be decomposed into recursive substructure:

$$= \sum_{y_{m-1}} (\exp s_m(y_m, y_{m-1})) \sum_{\mathbf{y}_{1:m-2}} \prod_{n=1}^{m-1} \exp s_n(y_n, y_{n-1})$$

$$= \sum_{y_{m-1}} (\exp s_m(y_m, y_{m-1})) \times \alpha_{m-1}(y_{m-1})$$

sum instead of max

Viterbi is a special case of the *max-product algorithm*,  
forward is a special case of the *sum-product algorithm*.

$$v_m(k) = \bigoplus_{k' \in \mathcal{Y}} s_m(k, k') \otimes v_{m-1}(k')$$

# Learning in CRFs

# Learning in CRFs

- As in logistic regression, gradient of the likelihood w.r.t. parameters is difference between observed and expected feature counts:

$$\frac{\delta \ell}{\delta \theta_j} = \sum_{i=1}^N E[f_j(\mathbf{w}^{(i)}, \mathbf{y})] - f_j(\mathbf{w}^{(i)}, \mathbf{y}^{(i)})$$

# Learning in CRFs

- As in logistic regression, gradient of the likelihood w.r.t. parameters is difference between observed and expected feature counts:

$$\frac{\delta \ell}{\delta \theta_j} = \sum_{i=1}^N E[f_j(\mathbf{w}^{(i)}, \mathbf{y})] - f_j(\mathbf{w}^{(i)}, \mathbf{y}^{(i)})$$

count of feature  $j$  for token sequence  $\mathbf{w}^{(i)}$ ,  
tag sequence  $\mathbf{y}^{(i)}$

# Learning in CRFs

- As in logistic regression, gradient of the likelihood w.r.t. parameters is difference between observed and expected feature counts:

$$\frac{\delta \ell}{\delta \theta_j} = \sum_{i=1}^N E[f_j(\mathbf{w}^{(i)}, \mathbf{y})] - f_j(\mathbf{w}^{(i)}, \mathbf{y}^{(i)})$$

count of feature  $j$  for token sequence  $\mathbf{w}^{(i)}$ ,  
tag sequence  $\mathbf{y}^{(i)}$

- Gradients can be computed by automatic differentiation!

# Learning in CRFs

- As in logistic regression, gradient of the likelihood w.r.t. parameters is difference between observed and expected feature counts:

$$\frac{\delta \ell}{\delta \theta_j} = \sum_{i=1}^N E[f_j(\mathbf{w}^{(i)}, \mathbf{y})] - f_j(\mathbf{w}^{(i)}, \mathbf{y}^{(i)})$$

count of feature  $j$  for token sequence  $\mathbf{w}^{(i)}$ ,  
tag sequence  $\mathbf{y}^{(i)}$

- Gradients can be computed by automatic differentiation!
- In the Olden Days, would use the **forward-backward algorithm** to compute expected counts.

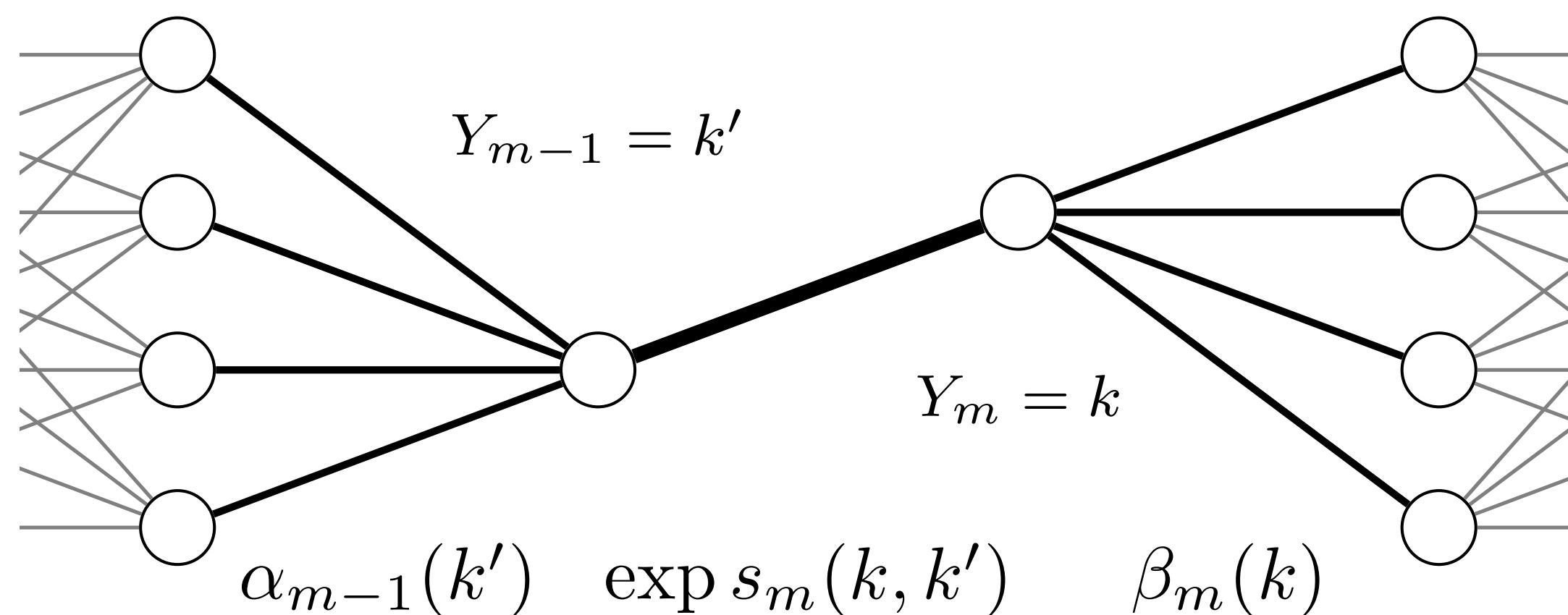
# Learning in CRFs

- As in logistic regression, gradient of the likelihood w.r.t. parameters is difference between observed and expected feature counts:

$$\frac{\delta \ell}{\delta \theta_j} = \sum_{i=1}^N E[f_j(\mathbf{w}^{(i)}, \mathbf{y})] - f_j(\mathbf{w}^{(i)}, \mathbf{y}^{(i)})$$

count of feature  $j$  for token sequence  $\mathbf{w}^{(i)}$ ,  
tag sequence  $\mathbf{y}^{(i)}$

- Gradients can be computed by automatic differentiation!
- In the Olden Days, would use the **forward-backward algorithm** to compute expected counts.



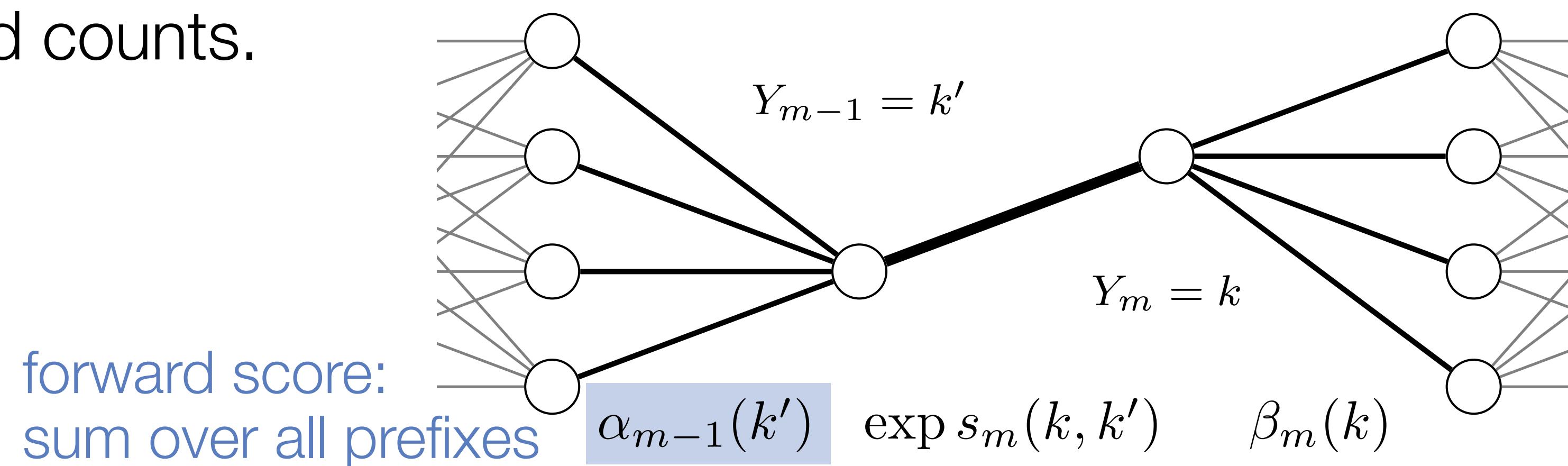
# Learning in CRFs

- As in logistic regression, gradient of the likelihood w.r.t. parameters is difference between observed and expected feature counts:

$$\frac{\delta \ell}{\delta \theta_j} = \sum_{i=1}^N E[f_j(\mathbf{w}^{(i)}, \mathbf{y})] - f_j(\mathbf{w}^{(i)}, \mathbf{y}^{(i)})$$

count of feature  $j$  for token sequence  $\mathbf{w}^{(i)}$ ,  
tag sequence  $\mathbf{y}^{(i)}$

- Gradients can be computed by automatic differentiation!
- In the Olden Days, would use the **forward-backward algorithm** to compute expected counts.



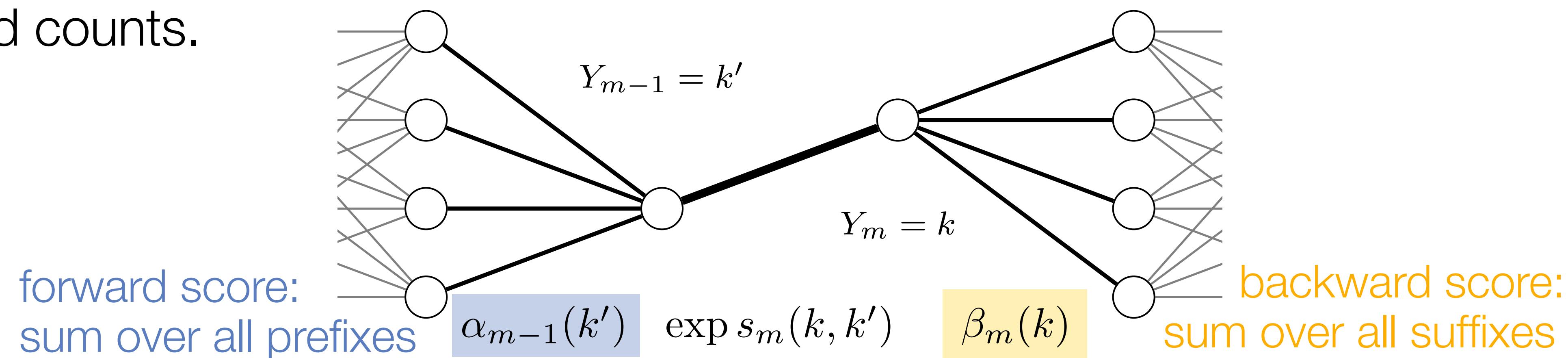
# Learning in CRFs

- As in logistic regression, gradient of the likelihood w.r.t. parameters is difference between observed and expected feature counts:

$$\frac{\delta \ell}{\delta \theta_j} = \sum_{i=1}^N E[f_j(\mathbf{w}^{(i)}, \mathbf{y})] - f_j(\mathbf{w}^{(i)}, \mathbf{y}^{(i)})$$

count of feature  $j$  for token sequence  $\mathbf{w}^{(i)}$ ,  
tag sequence  $\mathbf{y}^{(i)}$

- Gradients can be computed by automatic differentiation!
- In the Olden Days, would use the **forward-backward algorithm** to compute expected counts.



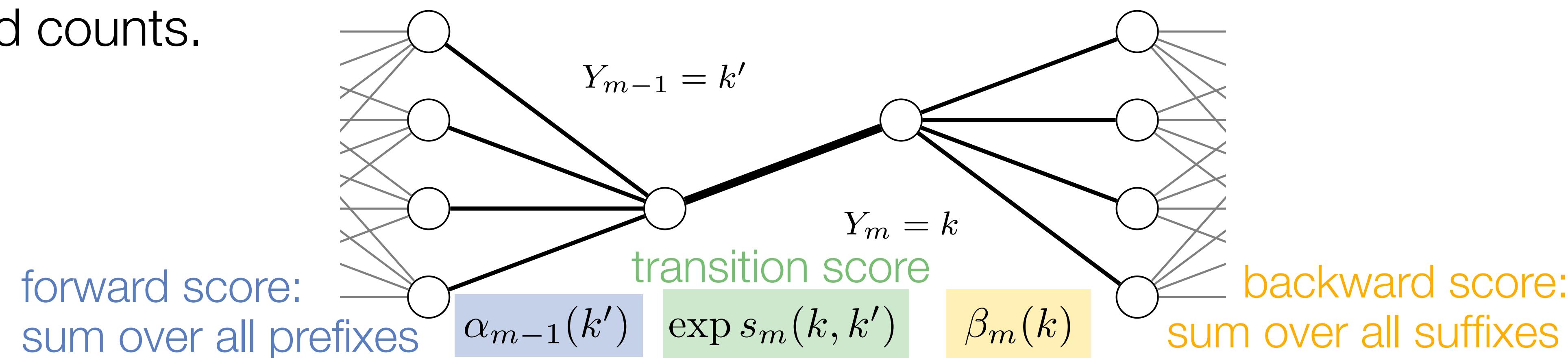
# Learning in CRFs

- As in logistic regression, gradient of the likelihood w.r.t. parameters is difference between observed and expected feature counts:

$$\frac{\delta \ell}{\delta \theta_j} = \sum_{i=1}^N E[f_j(\mathbf{w}^{(i)}, \mathbf{y})] - f_j(\mathbf{w}^{(i)}, \mathbf{y}^{(i)})$$

count of feature  $j$  for token sequence  $\mathbf{w}^{(i)}$ ,  
tag sequence  $\mathbf{y}^{(i)}$

- Gradients can be computed by automatic differentiation!
- In the Olden Days, would use the **forward-backward algorithm** to compute expected counts.



# Better features for sequence labeling?

- Until now: hand-engineered features:

# Better features for sequence labeling?

## ■ Until now: hand-engineered features:

Lexical	$f_{i \pm \{0,1,2,3\}}$ , $(m_{i-2,i-1})$ , $(m_{i-1,i})$ , $(m_{i-1,i+1})$ , $(m_{i,i+1})$ , $(m_{i+1,i+2})$ , $(m_{i-2,i-1,i})$ , $(m_{i-1,i,i+1})$ , $(m_{i,i+1,i+2})$ , $(m_{i-2,i-1,i+1})$ , $(m_{i-1,i+1,i+2})$
POS	$p_{i-\{3,2,1\}}$ , $a_{i+\{0,1,2,3\}}$ , $(p_{i-2,i-1})$ , $(a_{i+1,i+2})$ , $(p_{i-1}, a_{i+1})$ , $(p_{i-2}, p_{i-1}, a_i)$ , $(p_{i-2}, p_{i-1}, a_{i+1})$ , $(p_{i-1}, a_i, a_{i+1})$ , $(p_{i-1}, a_{i+1}, a_{i+2})$
Affix	$c_{:1}$ , $c_{:2}$ , $c_{:3}$ , $c_{n:}$ , $c_{n-1:}$ , $c_{n-2:}$ , $c_{n-3:}$
Binary	initial uppercase, all uppercase/lowercase, contains 1/2+ capital(s) not at the beginning, contains a (period/number/hyphen)

$w_i$  contains a particular prefix (from all prefixes of length  $\leq 4$ )

$w_i$  contains a particular suffix (from all suffixes of length  $\leq 4$ )

$w_i$  contains a number

$w_i$  contains an upper-case letter

$w_i$  contains a hyphen

$w_i$  is all upper case

$w_i$ 's word shape

$w_i$ 's short word shape

$w_i$  is upper case and has a digit and a dash (like *CFC-12*)

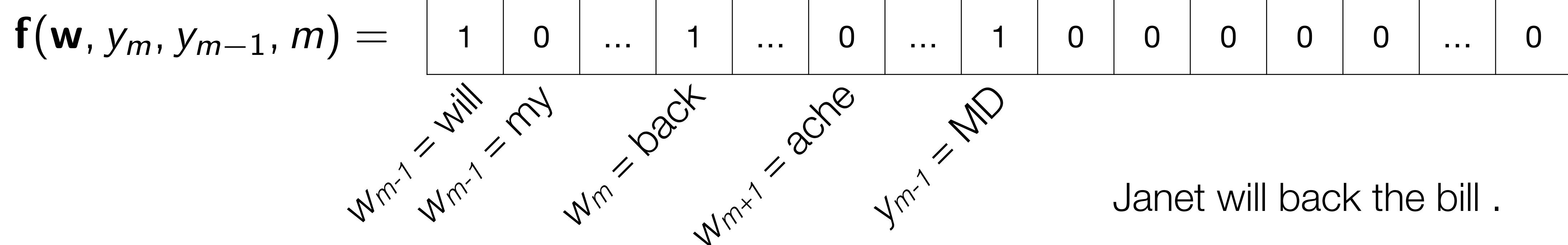
$w_i$  is upper case and followed within 3 words by Co., Inc., etc.

# Better features for sequence labeling?

- Until now: hand-engineered features:

Lexical	$f_{i \pm \{0,1,2,3\}}, (m_{i-2,i-1}), (m_{i-1,i}), (m_{i-1,i+1}), (m_{i,i+1}), (m_{i+1,i+2}), (m_{i-2,i-1,i}), (m_{i-1,i,i+1}), (m_{i,i+1,i+2}), (m_{i-2,i-1,i+1}), (m_{i-1,i+1,i+2})$
POS	$p_{i-\{3,2,1\}}, a_{i+\{0,1,2,3\}}, (p_{i-2,i-1}), (a_{i+1,i+2}), (p_{i-1}, a_{i+1}), (p_{i-2}, p_{i-1}, a_i), (p_{i-2}, p_{i-1}, a_{i+1}), (p_{i-1}, a_i, a_{i+1}), (p_{i-1}, a_{i+1}, a_{i+2})$
Affix	$c_{:1}, c_{:2}, c_{:3}, c_{n:,}, c_{n-1:,}, c_{n-2:,}, c_{n-3:}$
Binary	initial uppercase, all uppercase/lowercase, contains 1/2+ capital(s) not at the beginning, contains a (period/number/hyphen)

$w_i$  contains a particular prefix (from all prefixes of length  $\leq 4$ )  
 $w_i$  contains a particular suffix (from all suffixes of length  $\leq 4$ )  
 $w_i$  contains a number  
 $w_i$  contains an upper-case letter  
 $w_i$  contains a hyphen  
 $w_i$  is all upper case  
 $w_i$ 's word shape  
 $w_i$ 's short word shape  
 $w_i$  is upper case and has a digit and a dash (like *CFC-12*)  
 $w_i$  is upper case and followed within 3 words by Co., Inc., etc.



# Better features for sequence labeling?

- Until now: hand-engineered features:

**pros:**

**cons:**

$$f(\mathbf{w}, y_m, y_{m-1}, m) = \begin{bmatrix} 1 & 0 & \dots & 1 & \dots & 0 & \dots & 1 & 0 & 0 & 0 & 0 & 0 & \dots & 0 \end{bmatrix}$$

$w_{m-1} = \text{will}$   
 $w_{m-1} = \text{my}$   
 $w_m = \text{back}$   
 $w_{m+1} = \text{ache}$   
 $y_{m-1} = \text{MD}$

Janet will back the bill .

# Better features for sequence labeling?

- Until now: hand-engineered features:

## pros:

- interpretable, explainable
- can generalize well
- fast training and inference
- channel domain knowledge

## cons:

$$f(\mathbf{w}, y_m, y_{m-1}, m) = \begin{bmatrix} 1 & 0 & \dots & 1 & \dots & 0 & \dots & 1 & 0 & 0 & 0 & 0 & 0 & \dots & 0 \end{bmatrix}$$

$w_{m-1} = \text{will}$   
 $w_{m-1} = \text{my}$   
 $w_m = \text{back}$   
 $w_{m+1} = \text{ache}$   
 $y_{m-1} = \text{MD}$   
Janet will back the bill .

# Better features for sequence labeling?

- Until now: hand-engineered features:

## pros:

- interpretable, explainable
- can generalize well
- fast training and inference
- channel domain knowledge

## cons:

- can be sparse/high variance
- lack of shared representations
- task-specific
- worse performance

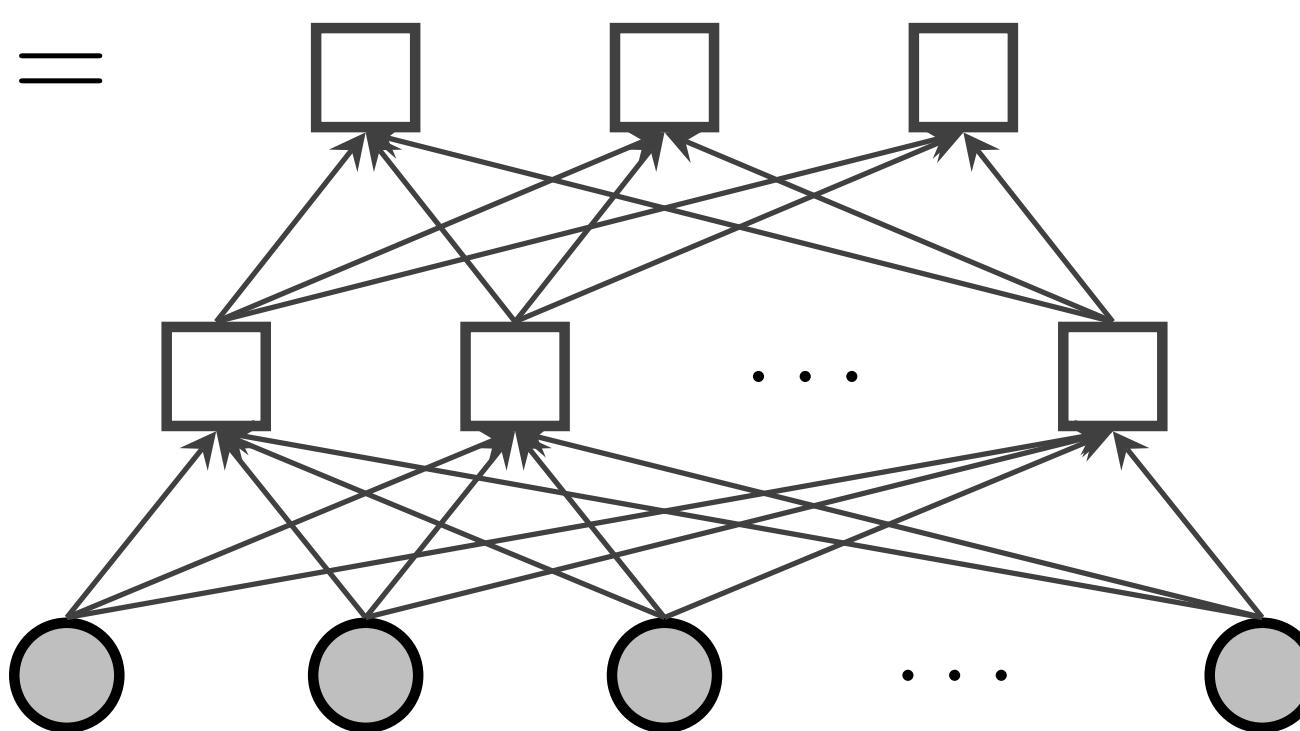
$$f(\mathbf{w}, y_m, y_{m-1}, m) = \begin{array}{cccccccccccccccccc} 1 & 0 & \dots & 1 & \dots & 0 & \dots & 1 & 0 & 0 & 0 & 0 & 0 & \dots & 0 \\ w_{m-1} = \text{will} & w_{m-1} = \text{my} & w_m = \text{back} & w_{m+1} = \text{ache} & y_{m-1} = \text{MD} & & & & & & & & & & & & & & \end{array}$$

Janet will back the bill .

# Neural sequence labeling

- Parameterize  $\mathbf{f}$  with a (deep) neural network.

$$\mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m) =$$



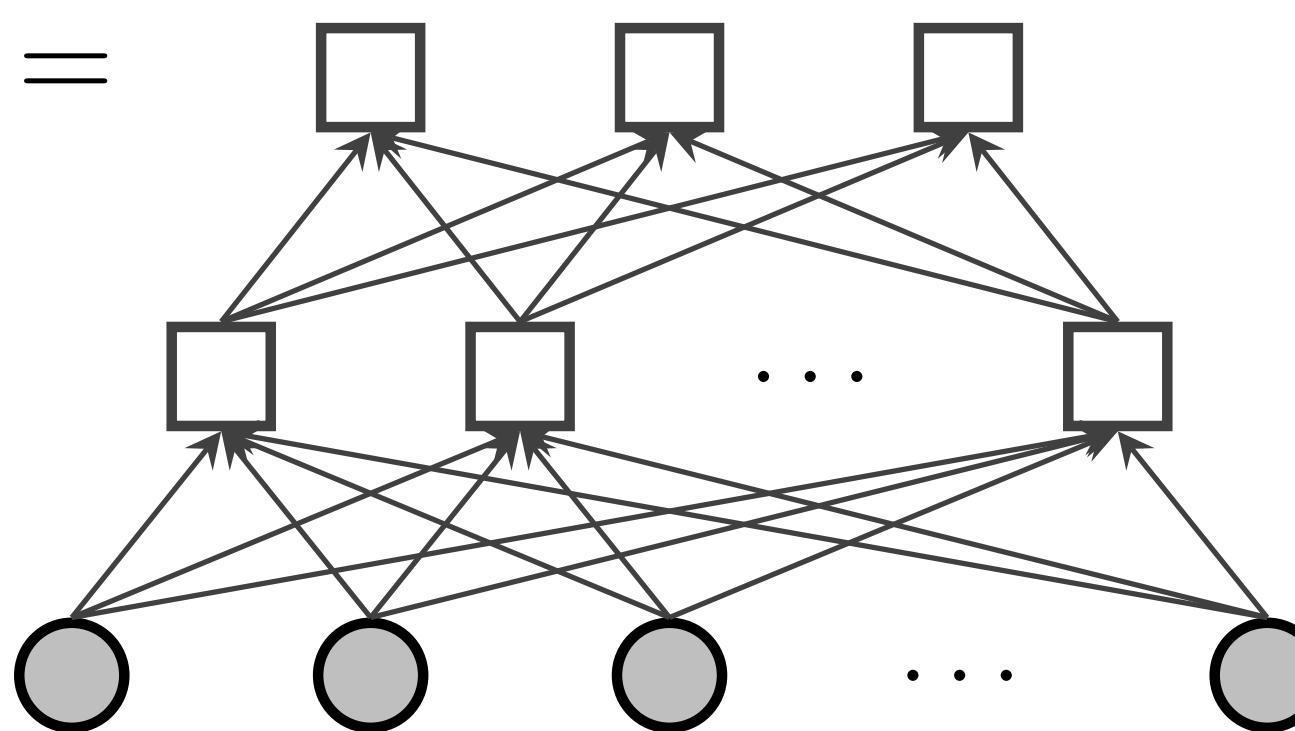
# Neural sequence labeling

- Parameterize  $\mathbf{f}$  with a (deep) neural network.

**pros:**

**cons:**

$$\mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m) =$$



# Neural sequence labeling

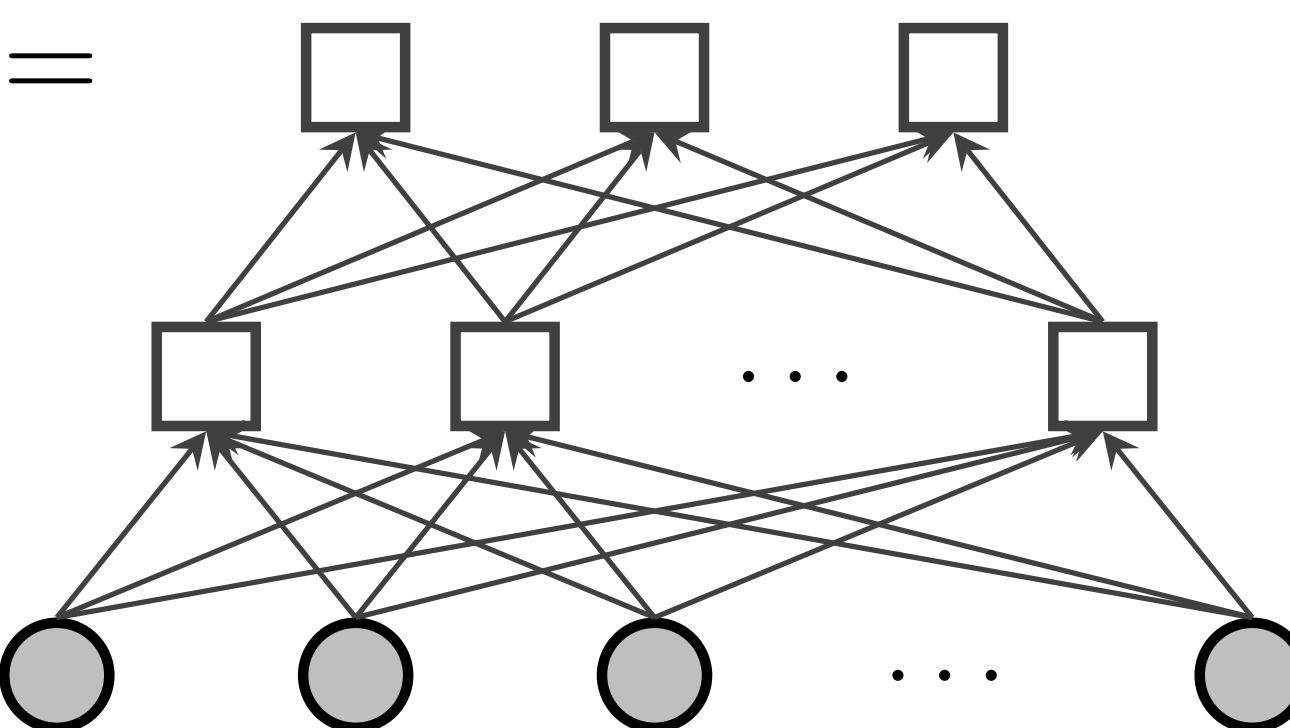
- Parameterize  $\mathbf{f}$  with a (deep) neural network.

## pros:

- shared representations
- channel external knowledge (e.g. word embeddings)
- high accuracy

## cons:

$$\mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m) =$$



# Neural sequence labeling

- Parameterize  $\mathbf{f}$  with a (deep) neural network.

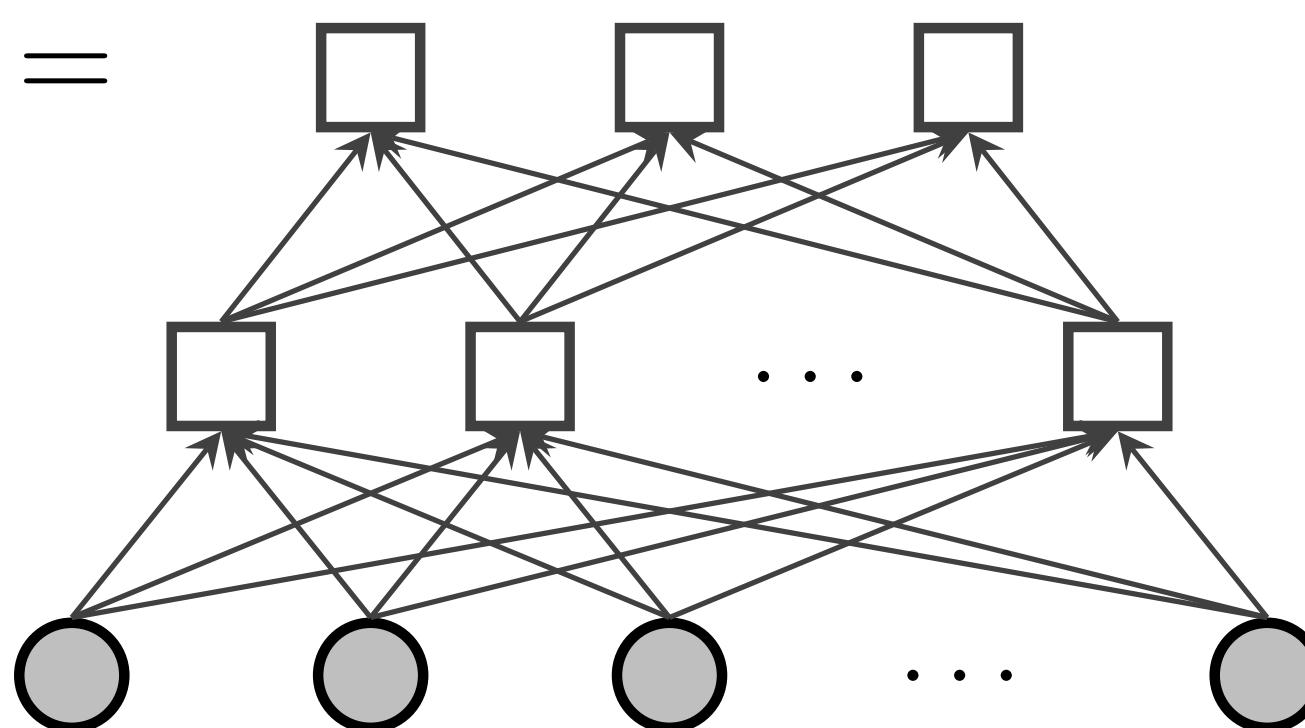
## pros:

- shared representations
- channel external knowledge (e.g. word embeddings)
- high accuracy

## cons:

- hard to interpret feature meaning, explain predictions
- optimization/hyperparameters
- prone to overfitting
- compute-heavy

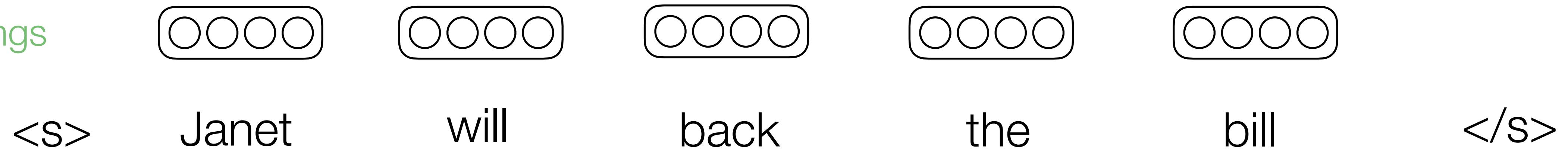
$$\mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m) =$$



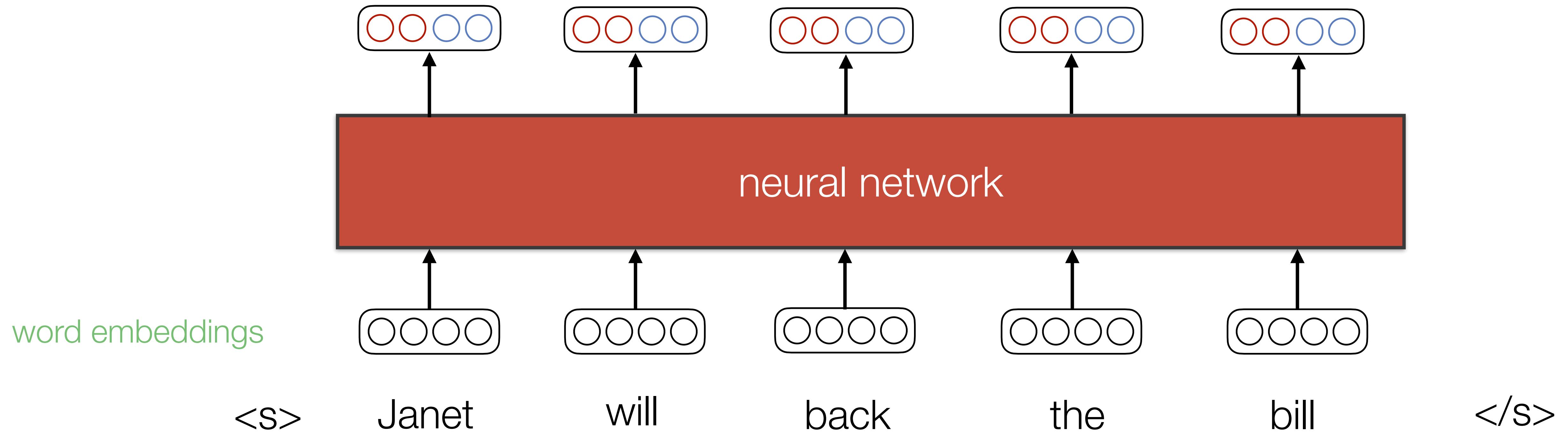
# Neural sequence labeling

# Neural sequence labeling

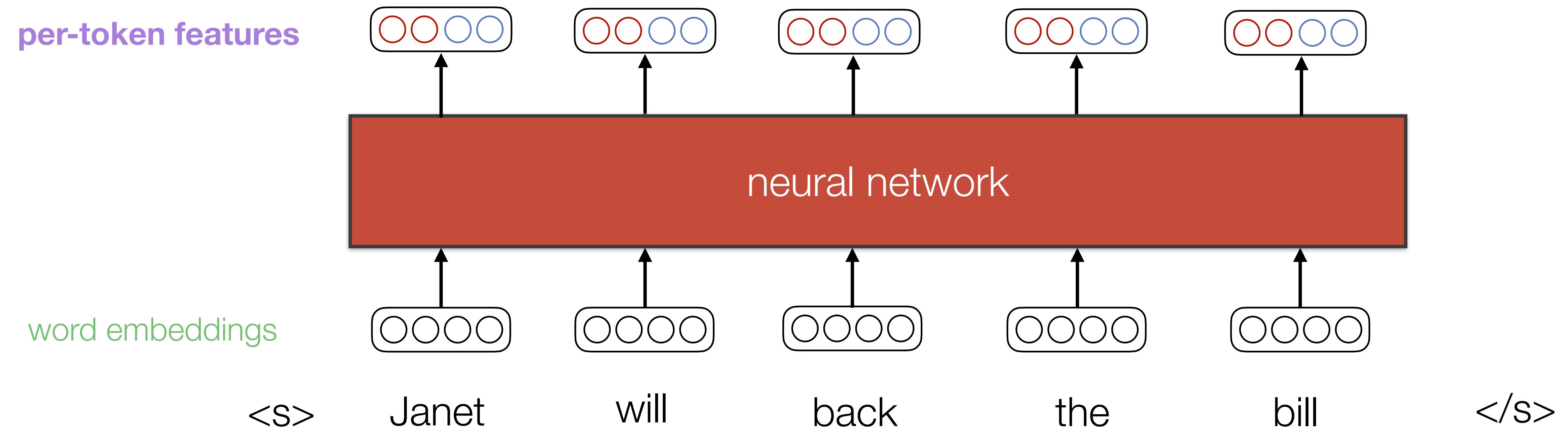
word embeddings



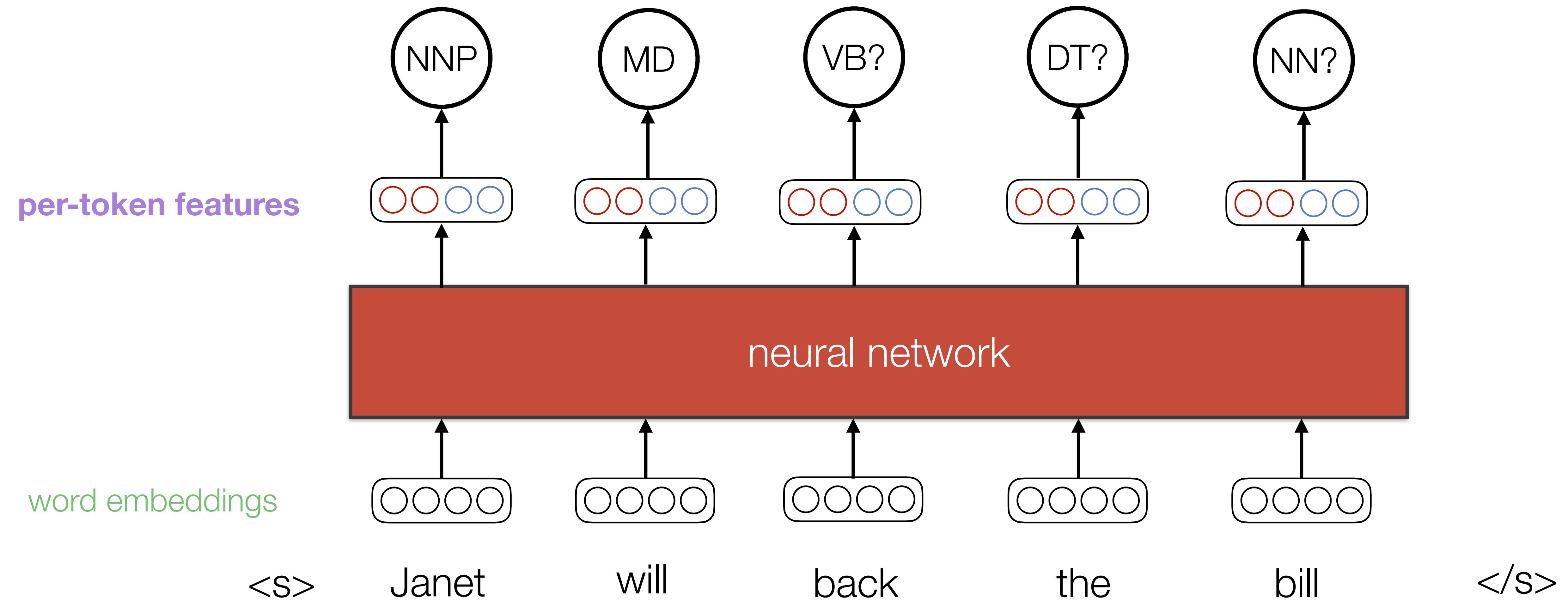
# Neural sequence labeling



# Neural sequence labeling

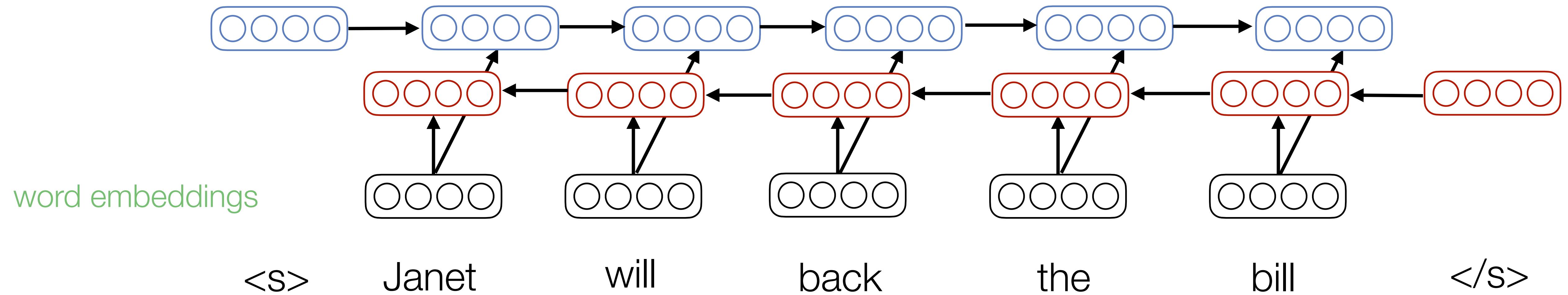


# Neural sequence labeling



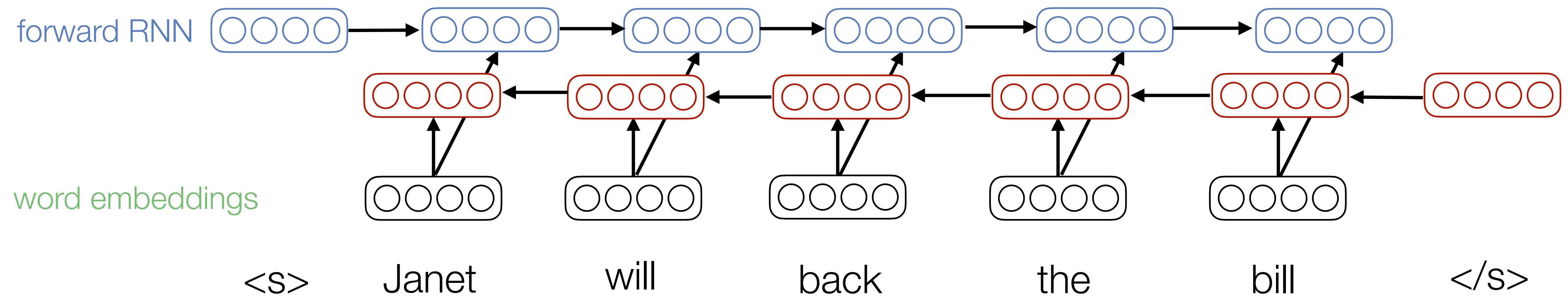
# Neural sequence labeling

## Bidirectional RNNs



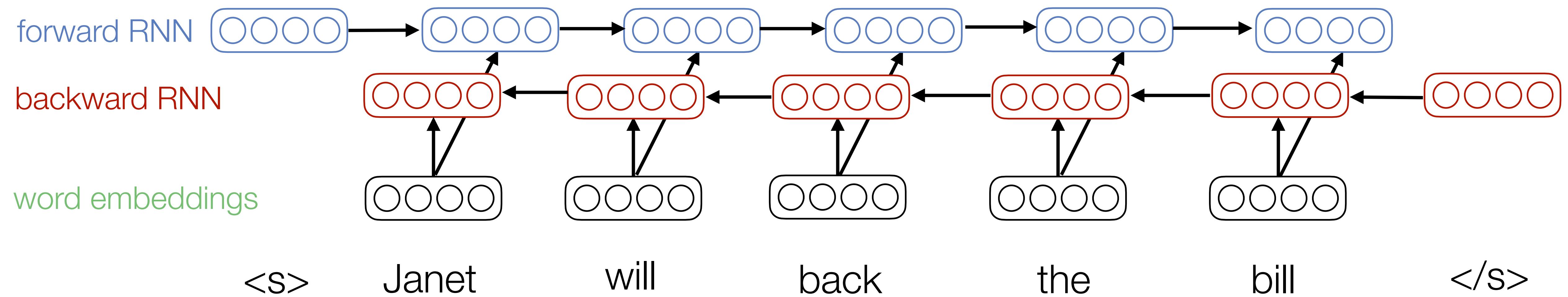
# Neural sequence labeling

## Bidirectional RNNs



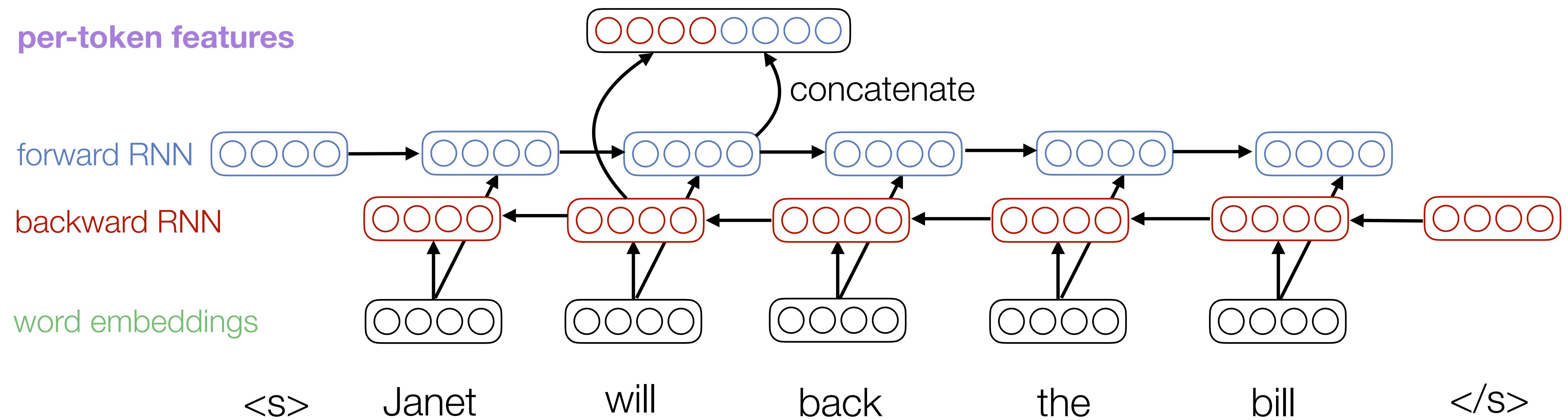
# Neural sequence labeling

## Bidirectional RNNs



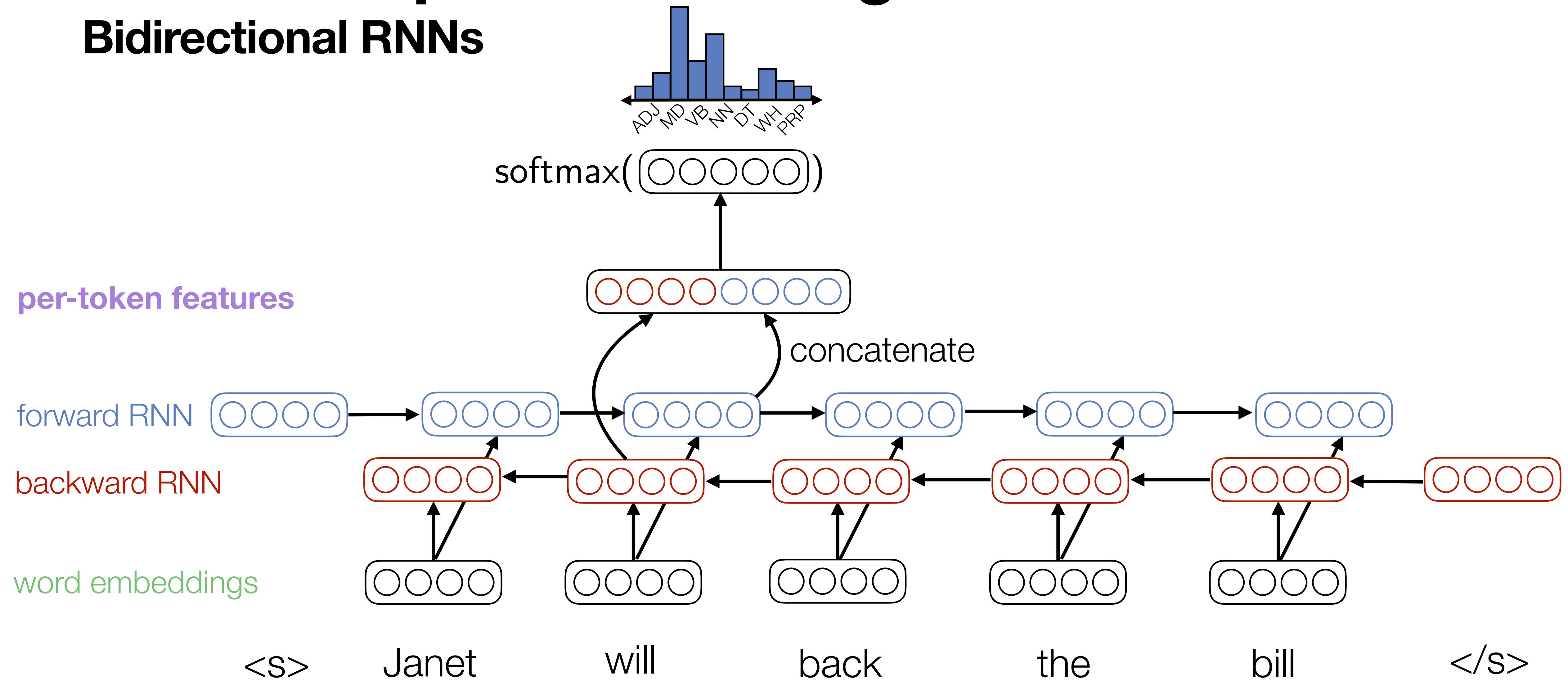
# Neural sequence labeling

## Bidirectional RNNs



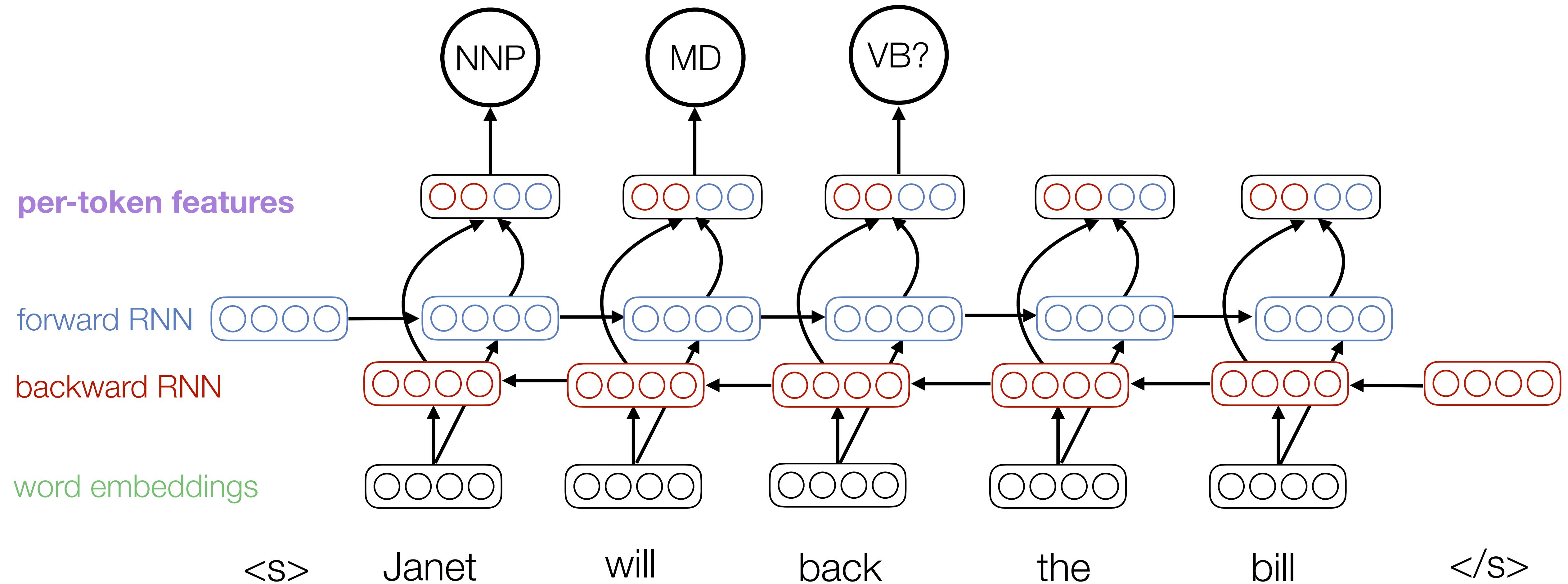
# Neural sequence labeling

## Bidirectional RNNs



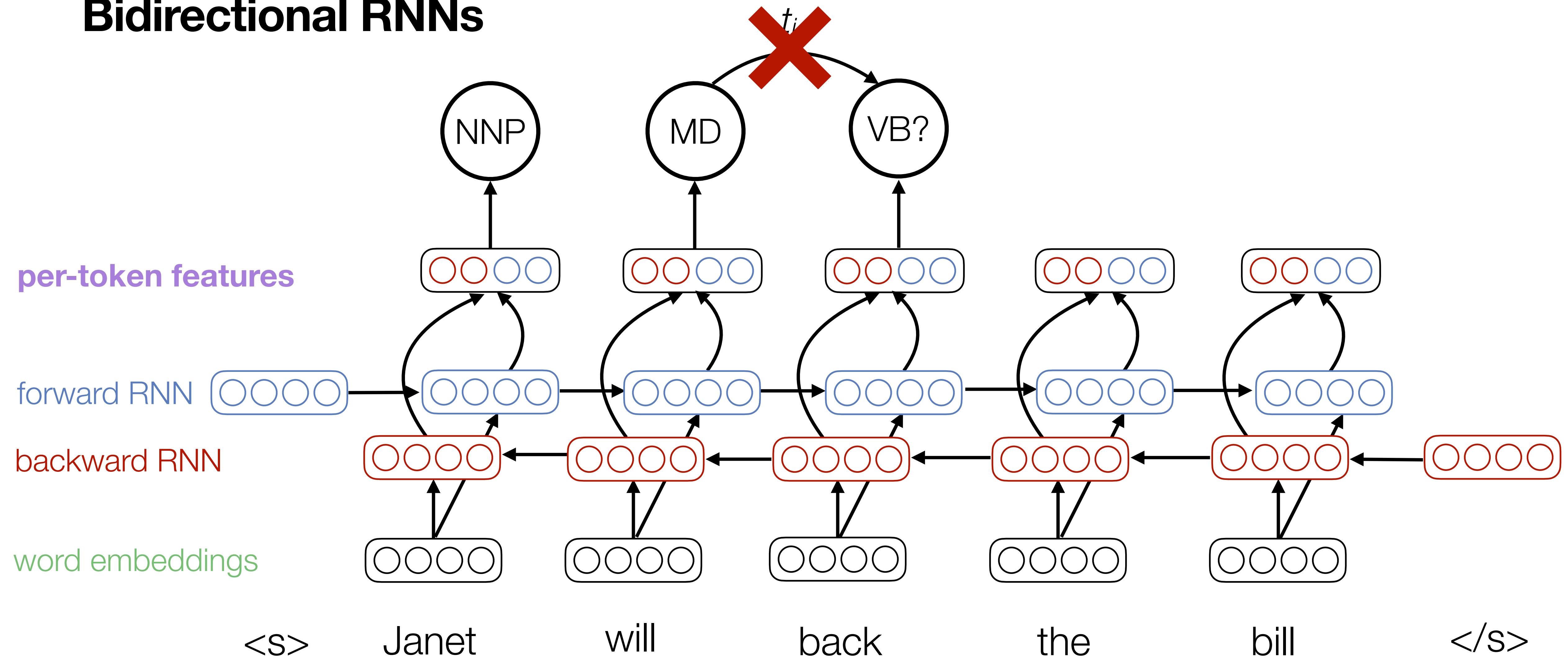
# Neural sequence labeling

## Bidirectional RNNs



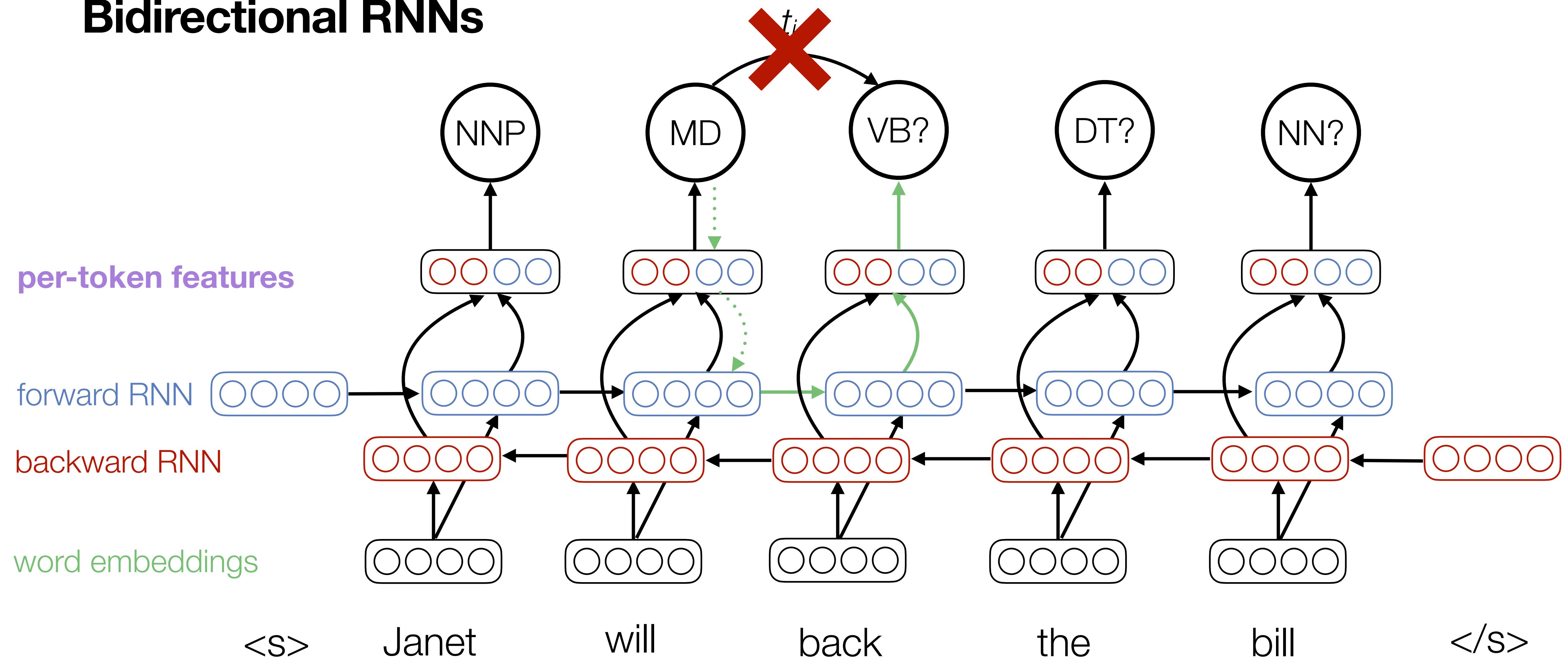
# Neural sequence labeling

## Bidirectional RNNs



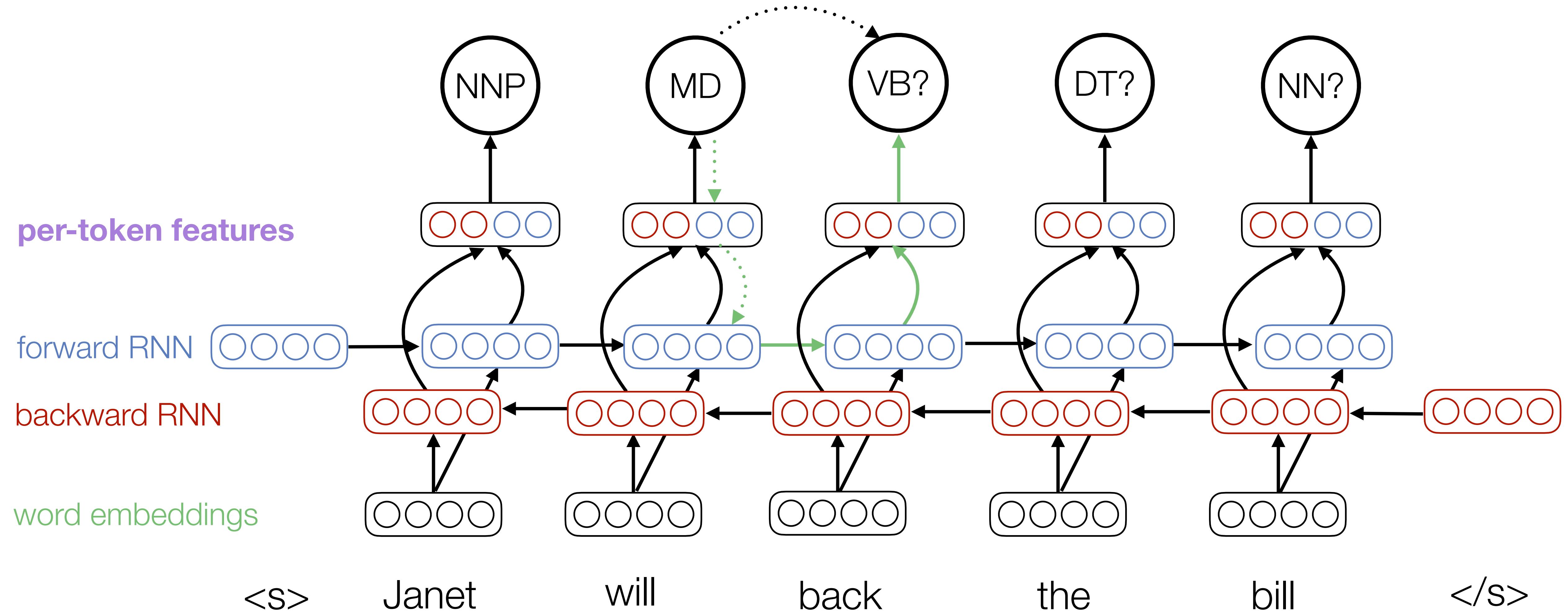
# Neural sequence labeling

## Bidirectional RNNs



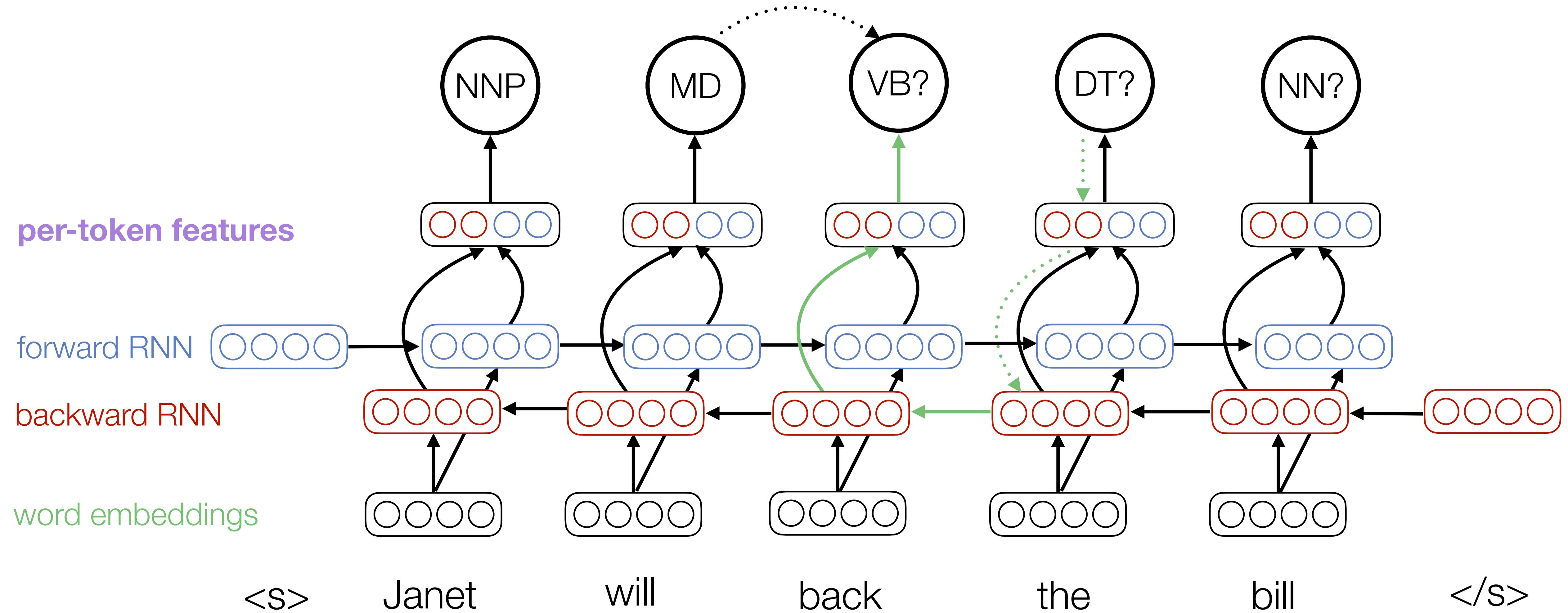
# Neural sequence labeling

## Bidirectional RNNs



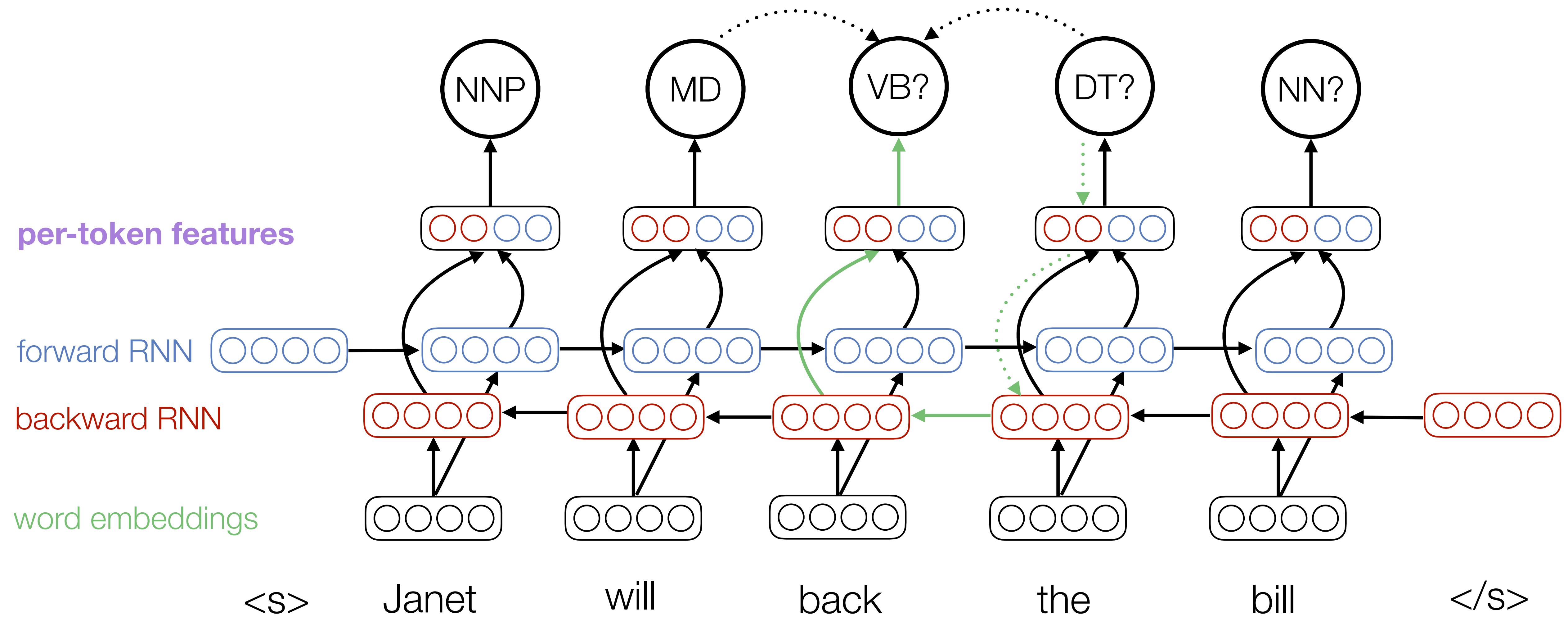
# Neural sequence labeling

## Bidirectional RNNs



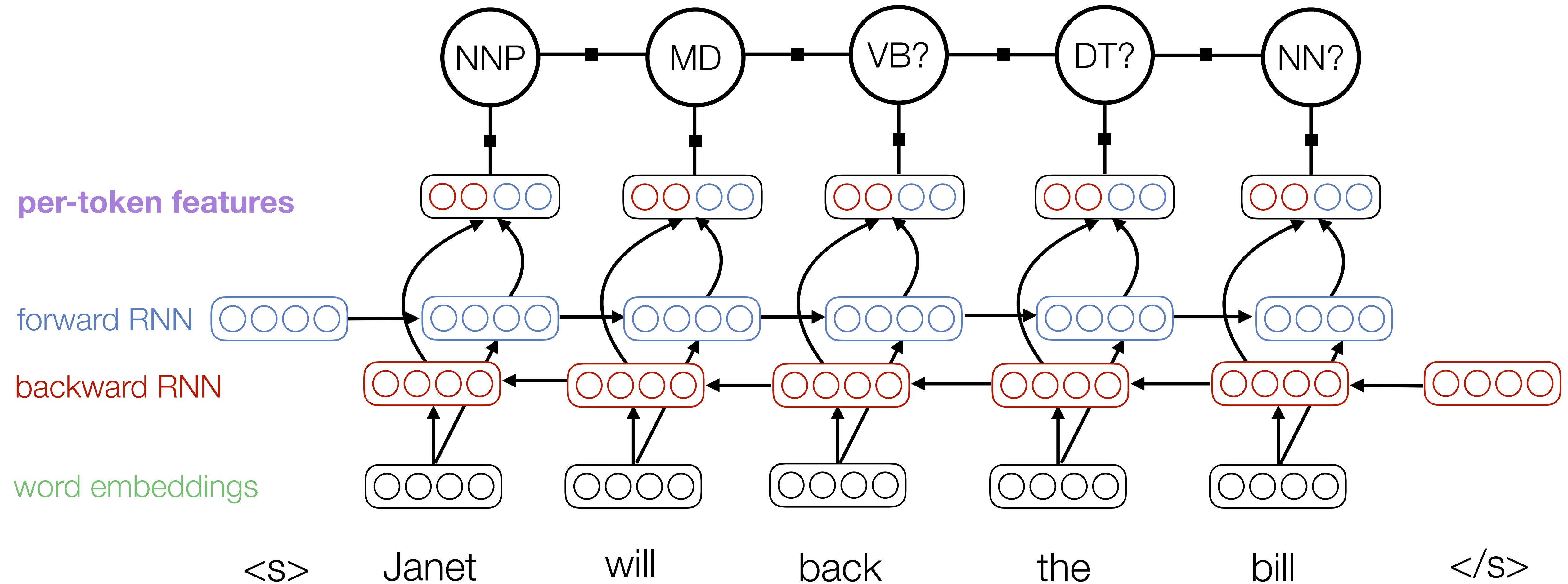
# Neural sequence labeling

## Bidirectional RNNs

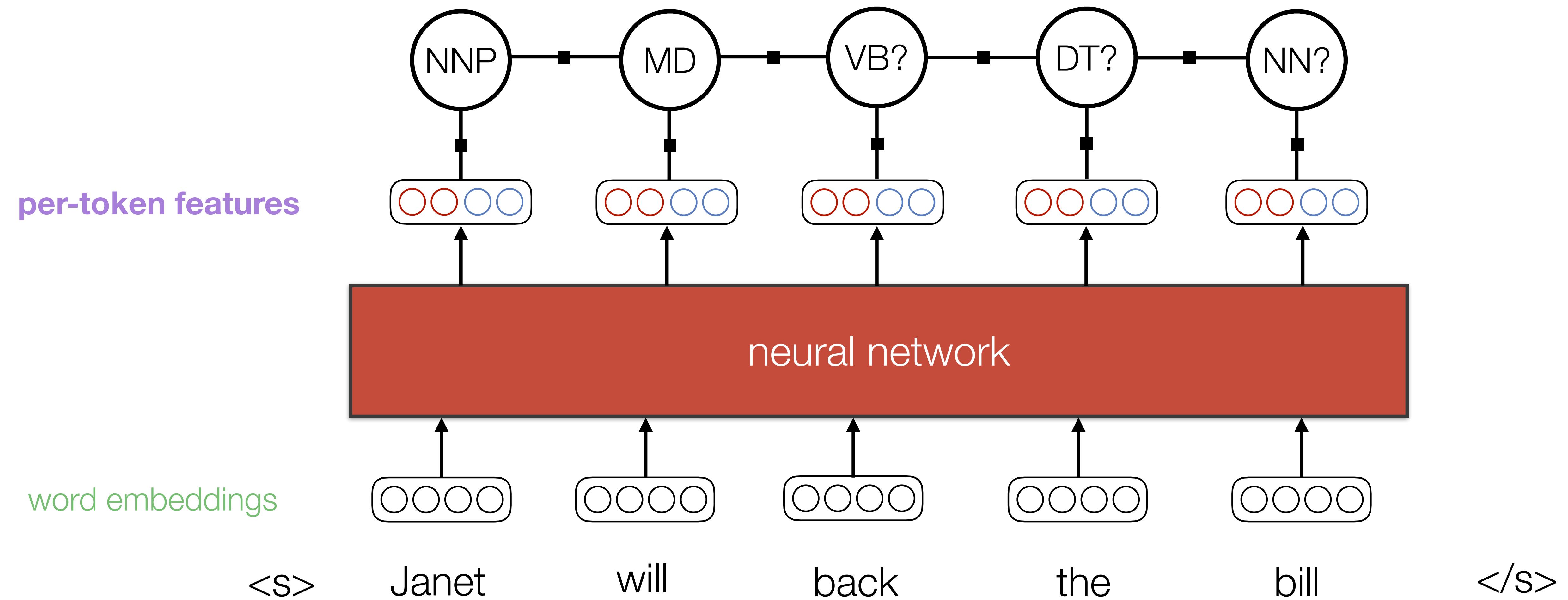


# Neural sequence labeling

## Bidirectional RNN-CRFs



# Neural sequence labeling



# Convolutional neural networks

# Convolutional neural networks

- In NLP, CNNs merge information across contiguous, fixed-width spans of tokens.

# Convolutional neural networks

- In NLP, CNNs merge information across contiguous, fixed-width spans of tokens.
- Unlike computer vision, in NLP we use 1D CNNs.

# Convolutional neural networks

- In NLP, CNNs merge information across contiguous, fixed-width spans of tokens.
- Unlike computer vision, in NLP we use 1D CNNs.
- For sentence/document classification: **pooling function** over representations.

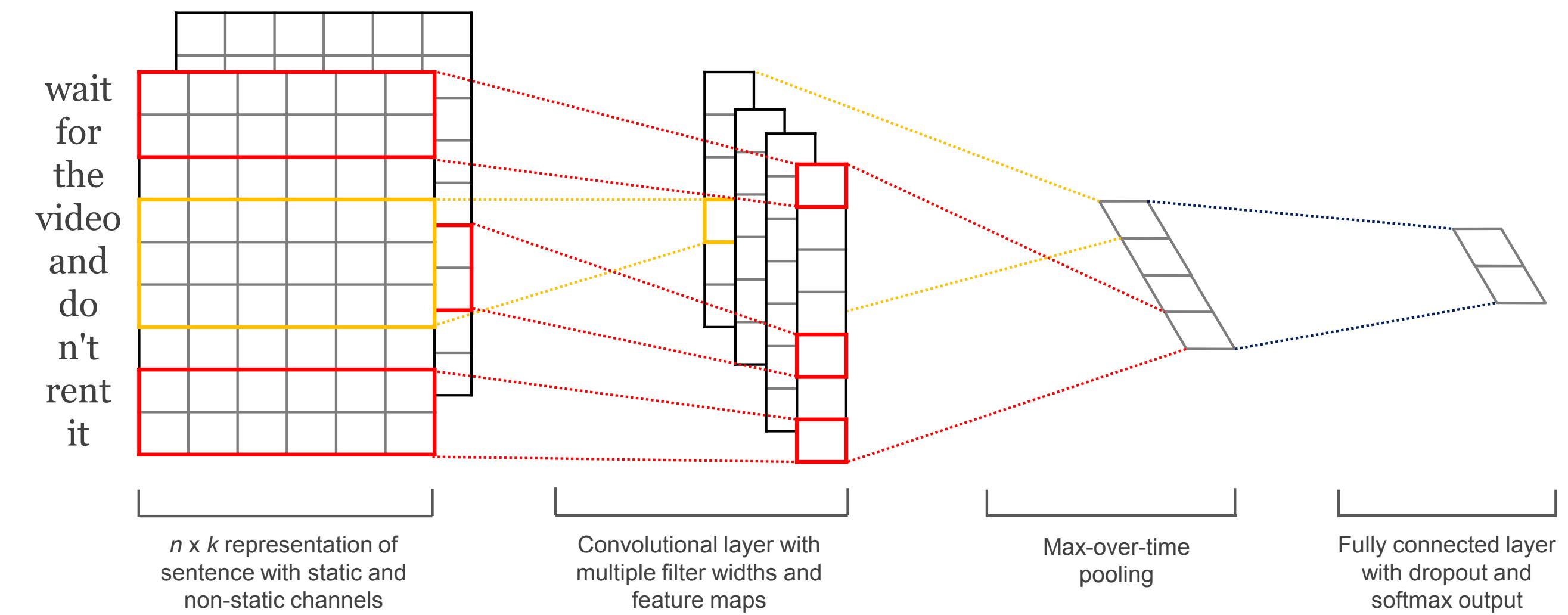


Figure from: [Yoon Kim. Convolutional Neural Networks for Sentence Classification. EMNLP 2014.](#)

# Convolutional neural networks

- In NLP, CNNs merge information across contiguous, fixed-width spans of tokens.
- Unlike computer vision, in NLP we use 1D CNNs.
- For sentence/document classification: **pooling function** over representations.
  - For example: sum, average. Most common: **max pooling** (over time).

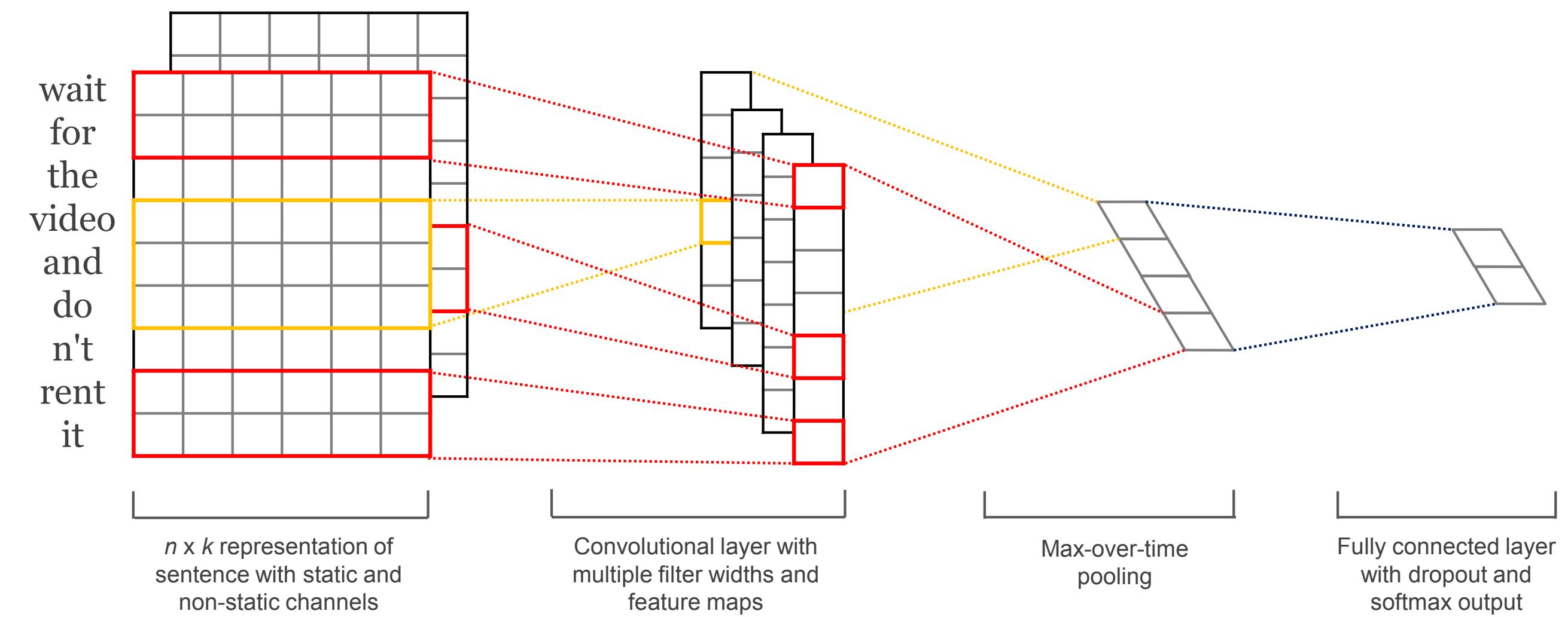
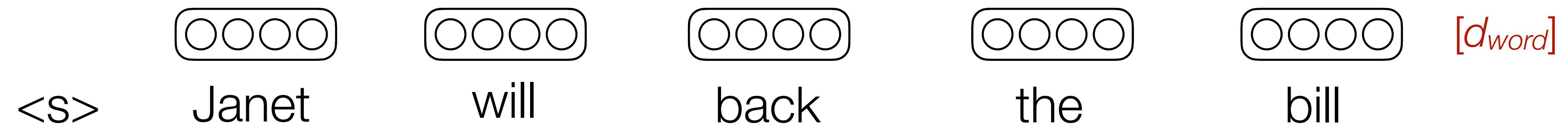


Figure from: [Yoon Kim. Convolutional Neural Networks for Sentence Classification. EMNLP 2014.](#)

# Convolutional neural networks

- In NLP, CNNs merge information across contiguous, fixed-width spans of tokens.

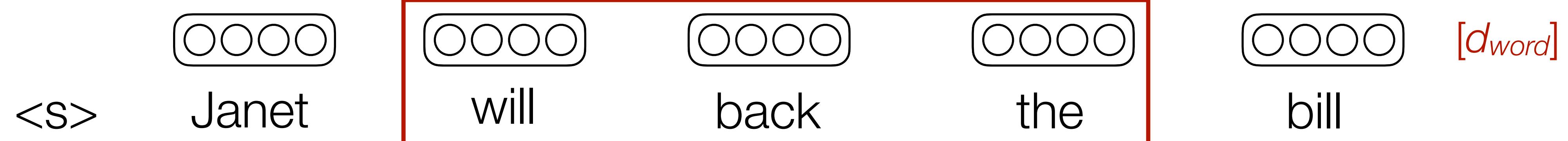
**dims:**



# Convolutional neural networks

- In NLP, CNNs merge information across contiguous, fixed-width spans of tokens.

**dims:**

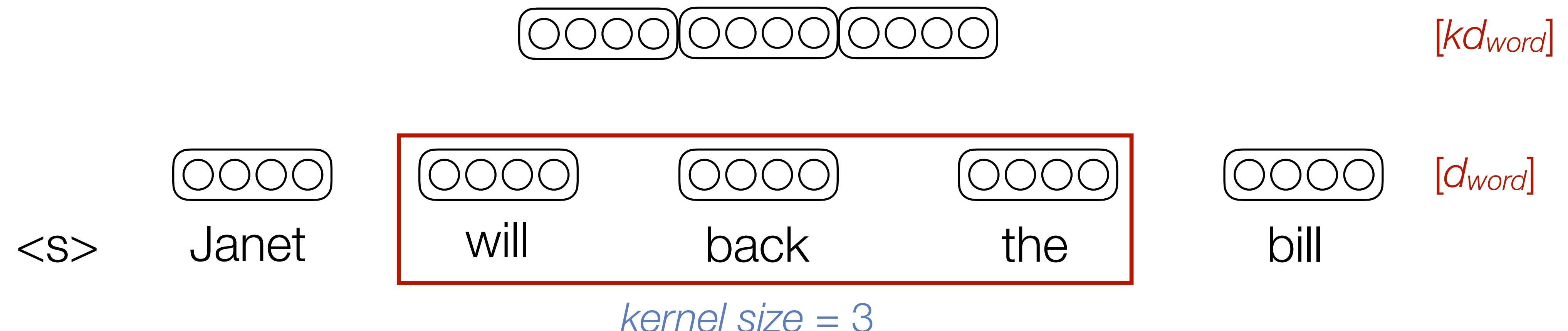


*kernel size = 3*

# Convolutional neural networks

- In NLP, CNNs merge information across contiguous, fixed-width spans of tokens.

**dims:**



# Convolutional neural networks

- In NLP, CNNs merge information across contiguous, fixed-width spans of tokens.

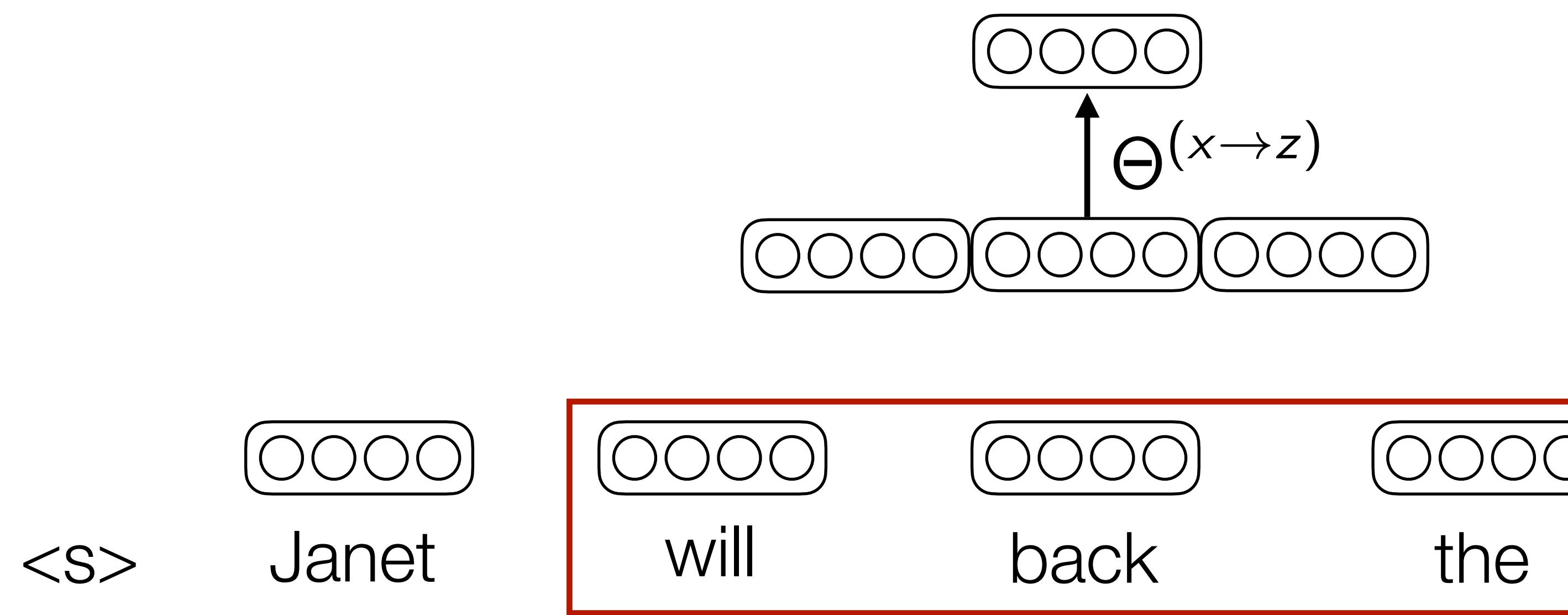
**dims:**

$[d_z]$  # filters

$[kd_{word} \times d_z]$

$[kd_{word}]$

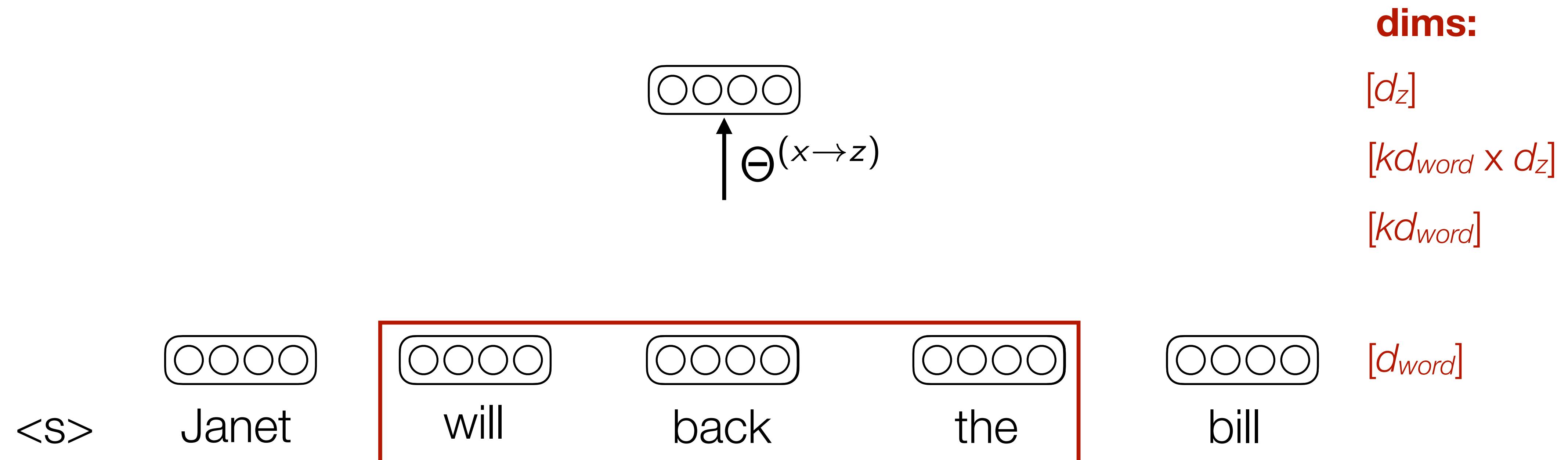
$[d_{word}]$



*kernel size = 3*

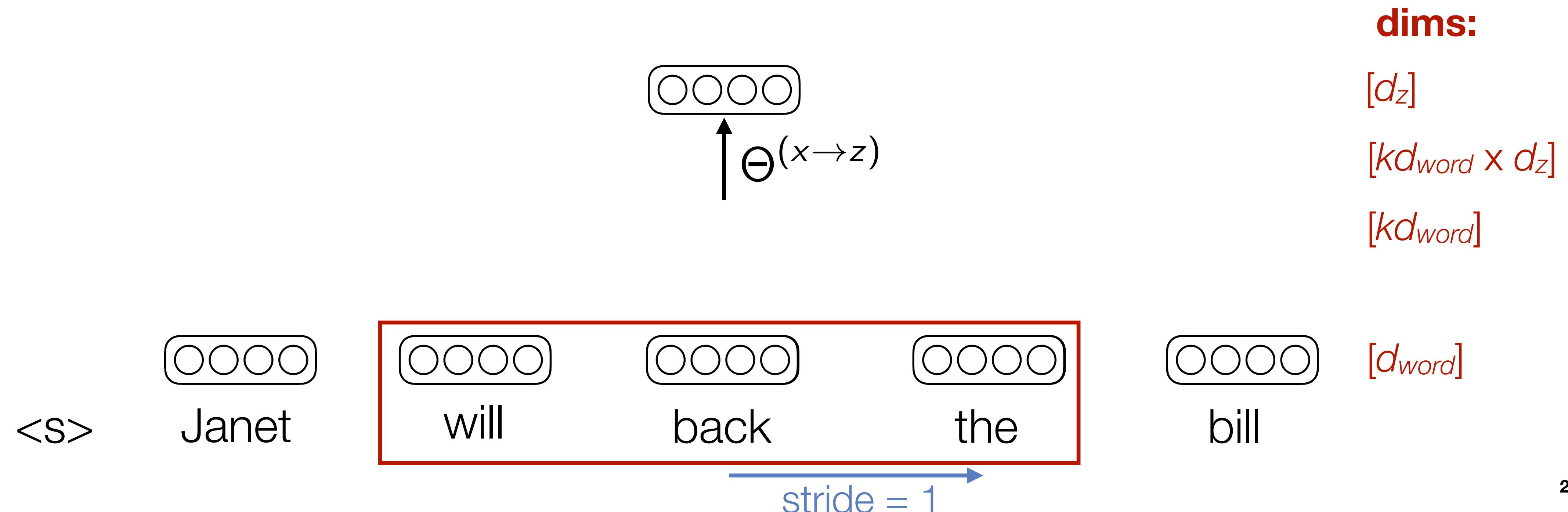
# Convolutional neural networks

- In NLP, CNNs merge information across contiguous, fixed-width spans of tokens.



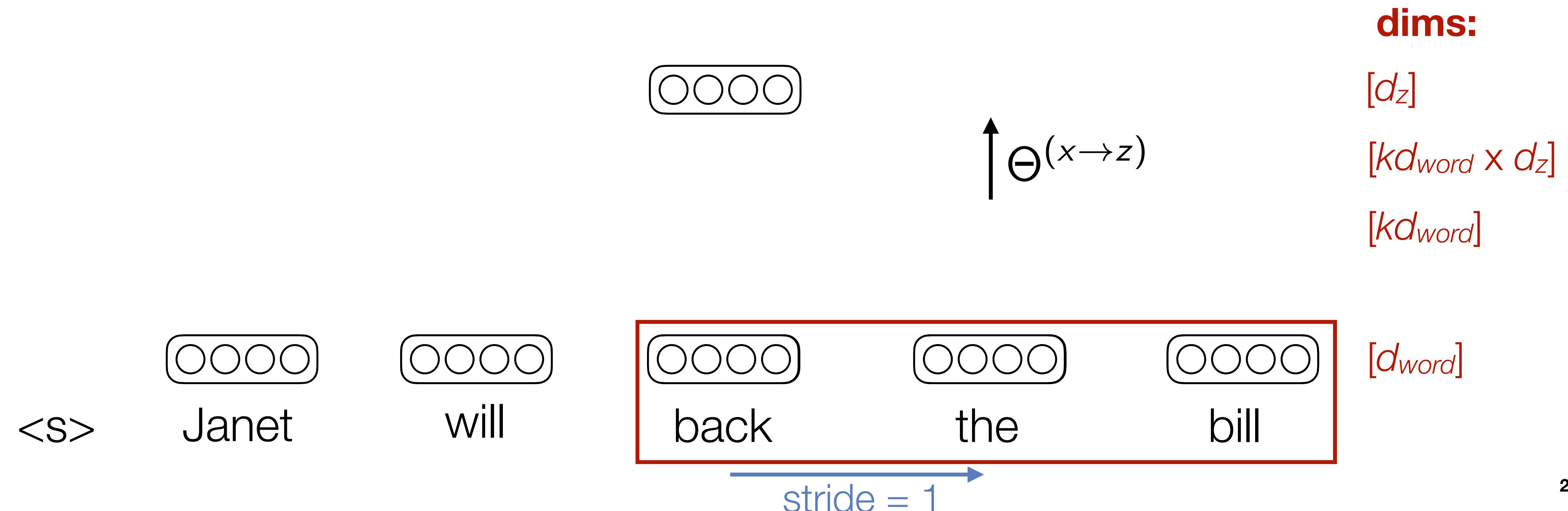
# Convolutional neural networks

- In NLP, CNNs merge information across contiguous, fixed-width spans of tokens.



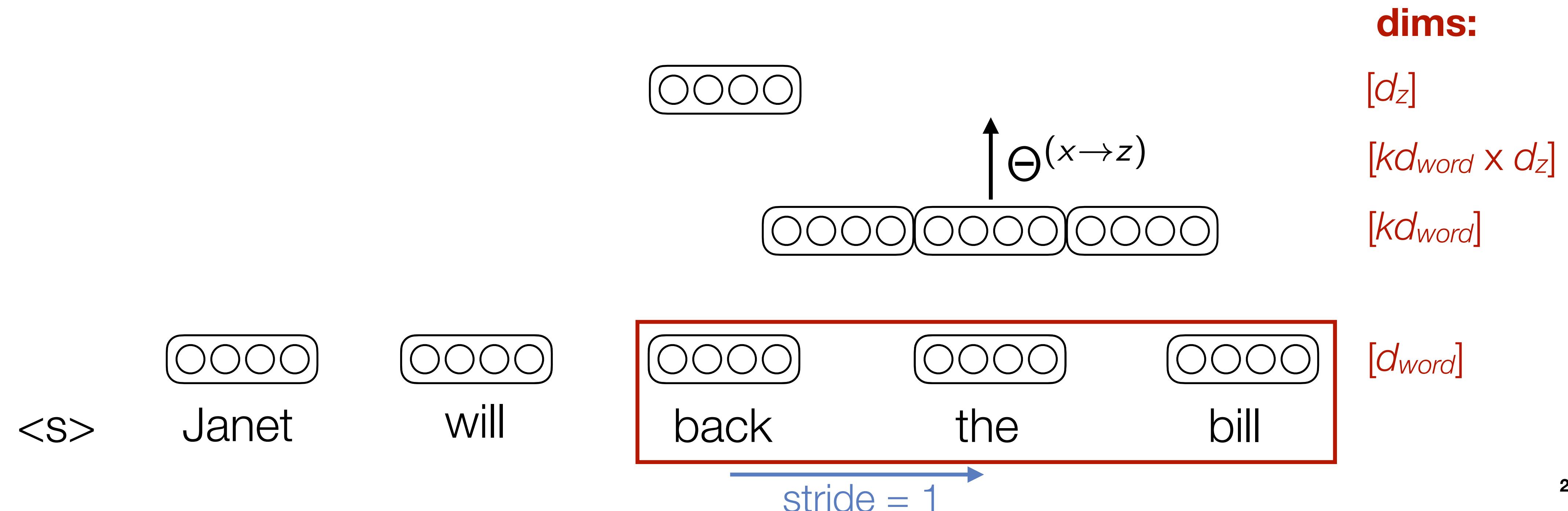
# Convolutional neural networks

- In NLP, CNNs merge information across contiguous, fixed-width spans of tokens.



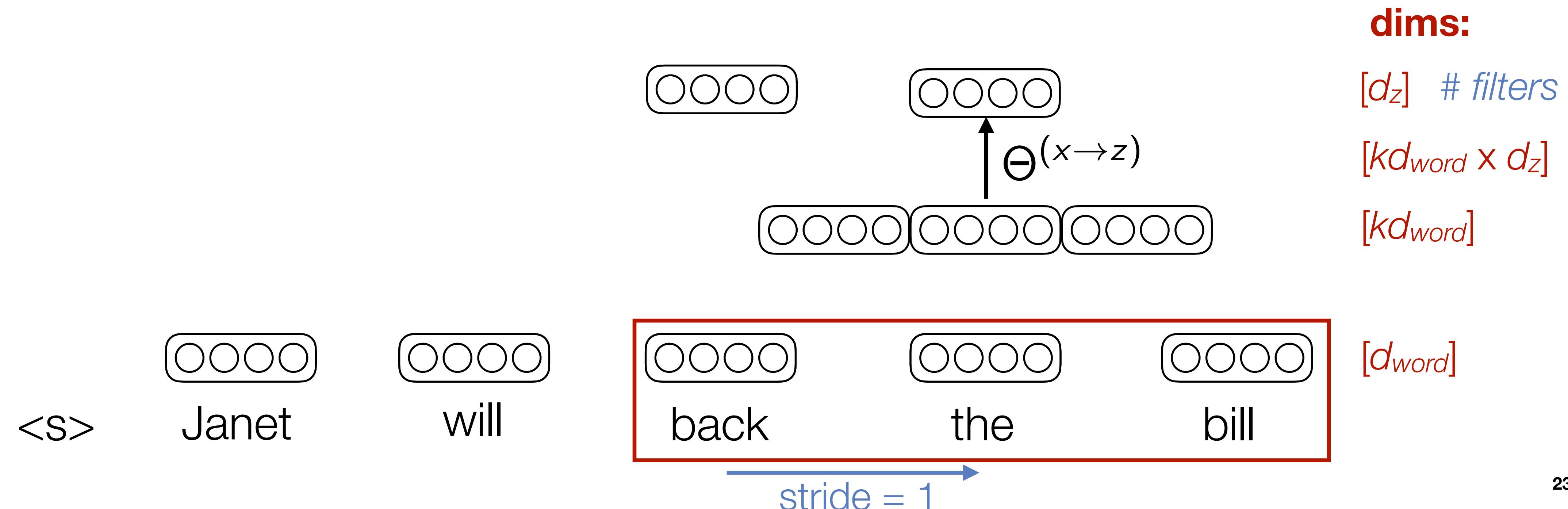
# Convolutional neural networks

- In NLP, CNNs merge information across contiguous, fixed-width spans of tokens.

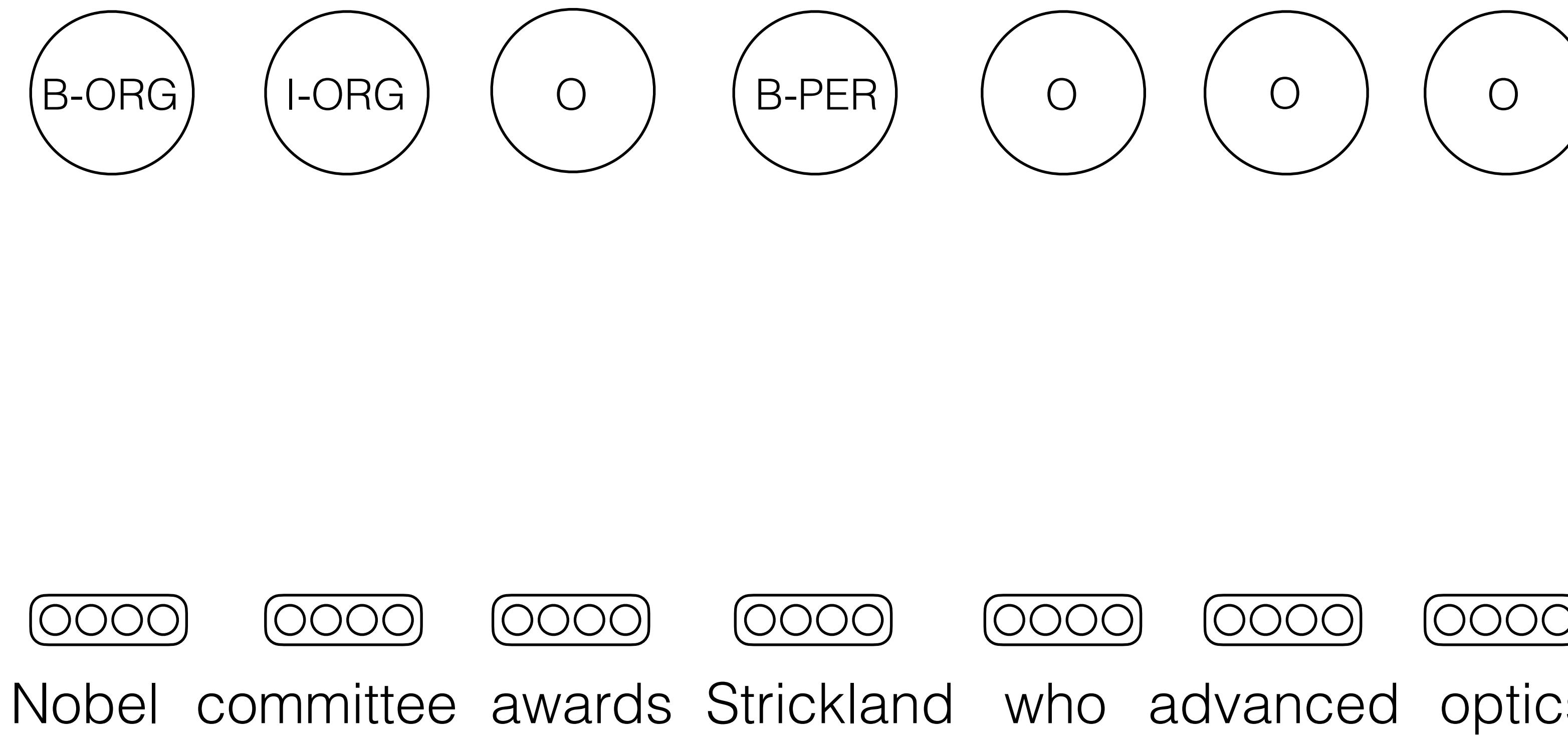


# Convolutional neural networks

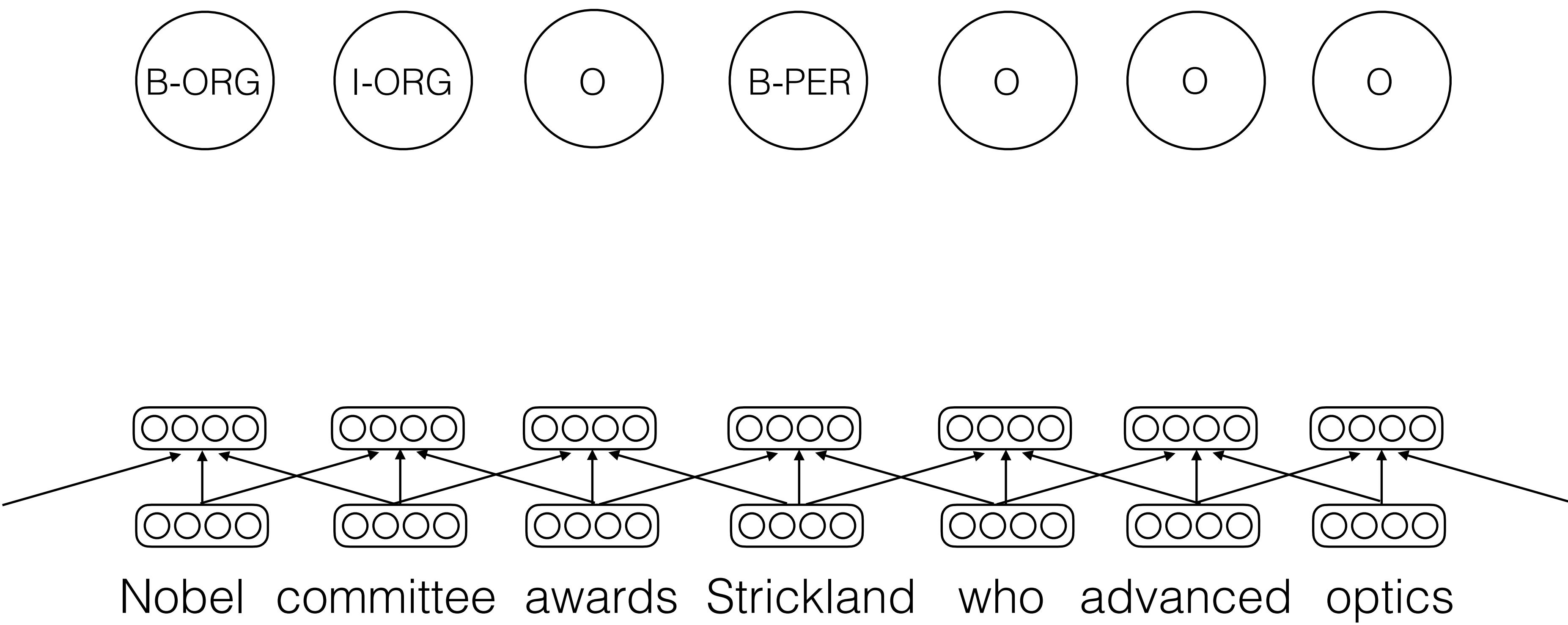
- In NLP, CNNs merge information across contiguous, fixed-width spans of tokens.



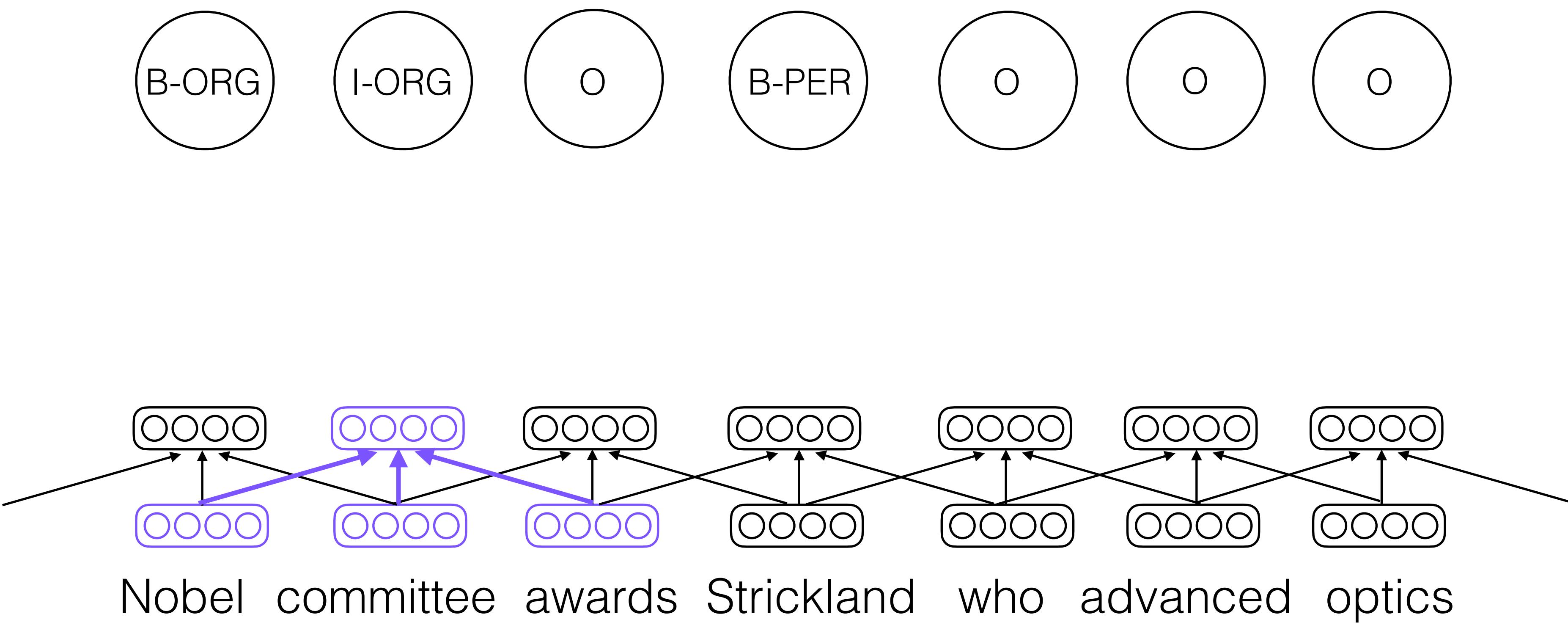
# Sequence labeling w/ CNNs



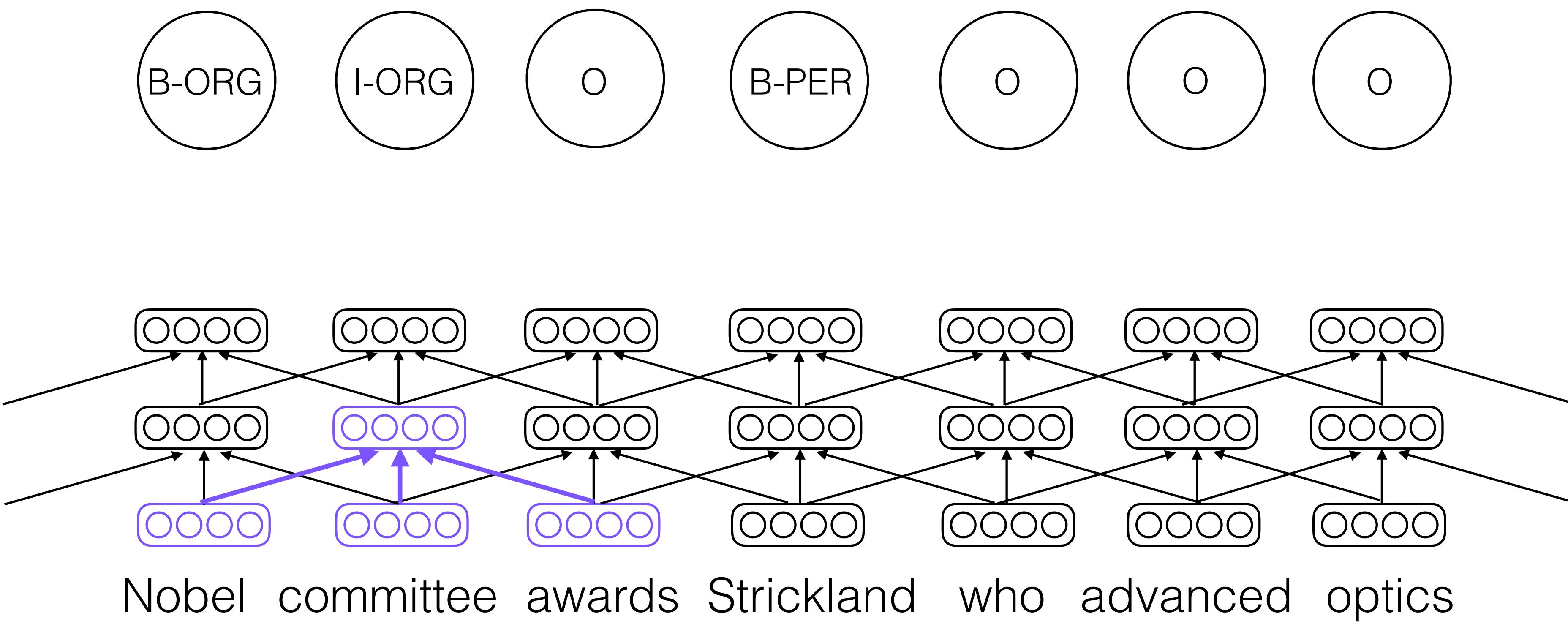
# Sequence labeling w/ CNNs



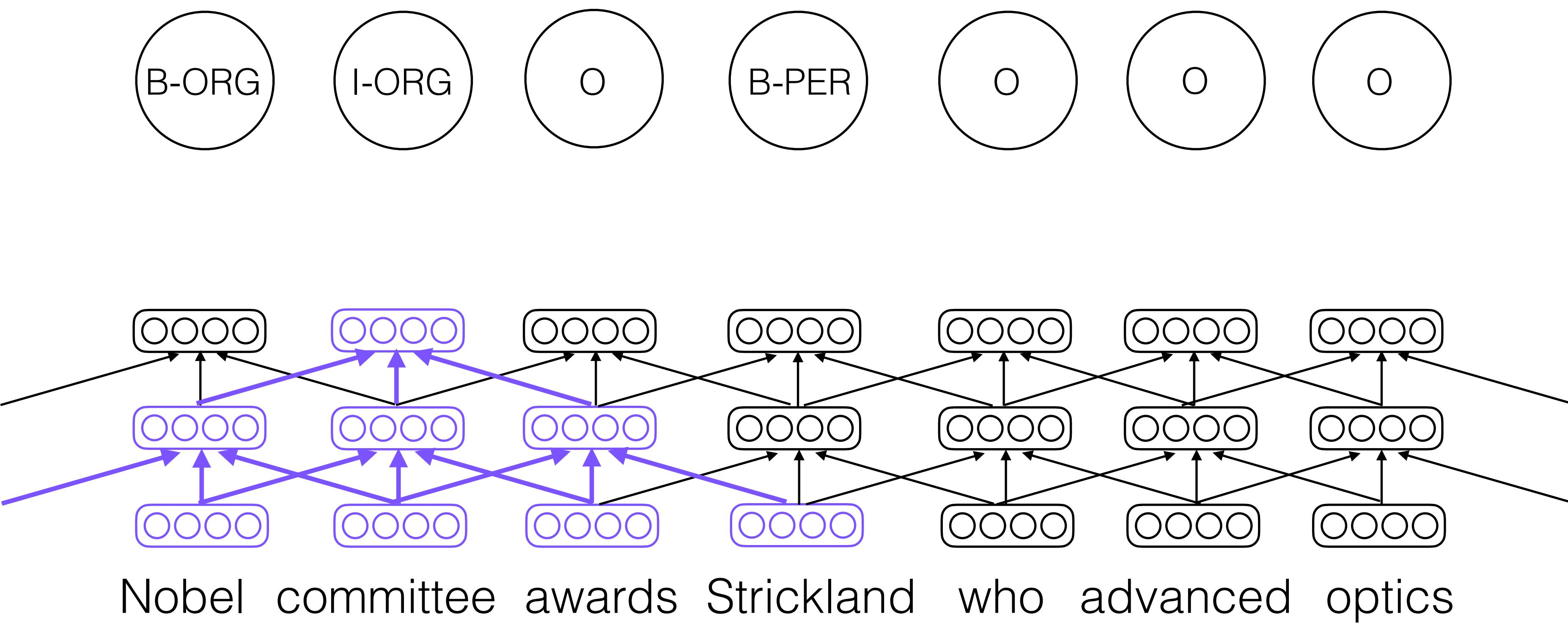
# Sequence labeling w/ CNNs



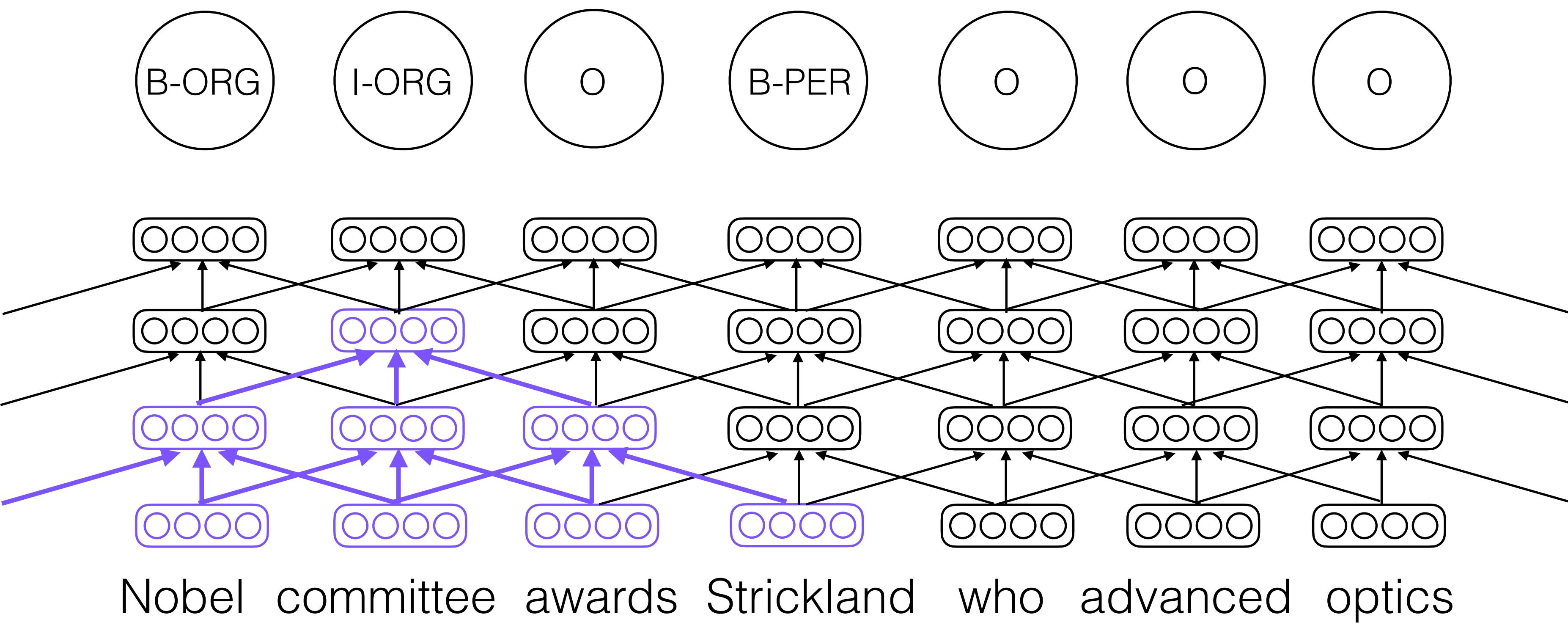
# Sequence labeling w/ CNNs



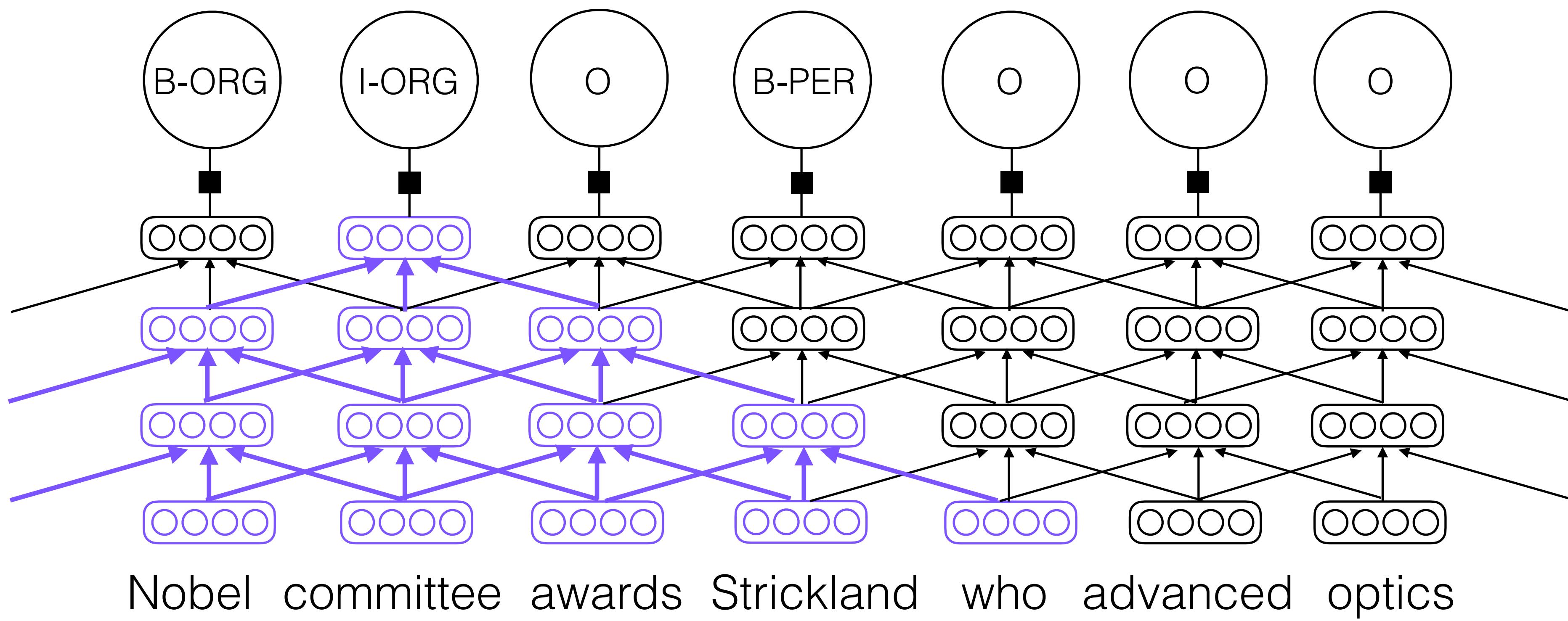
# Sequence labeling w/ CNNs



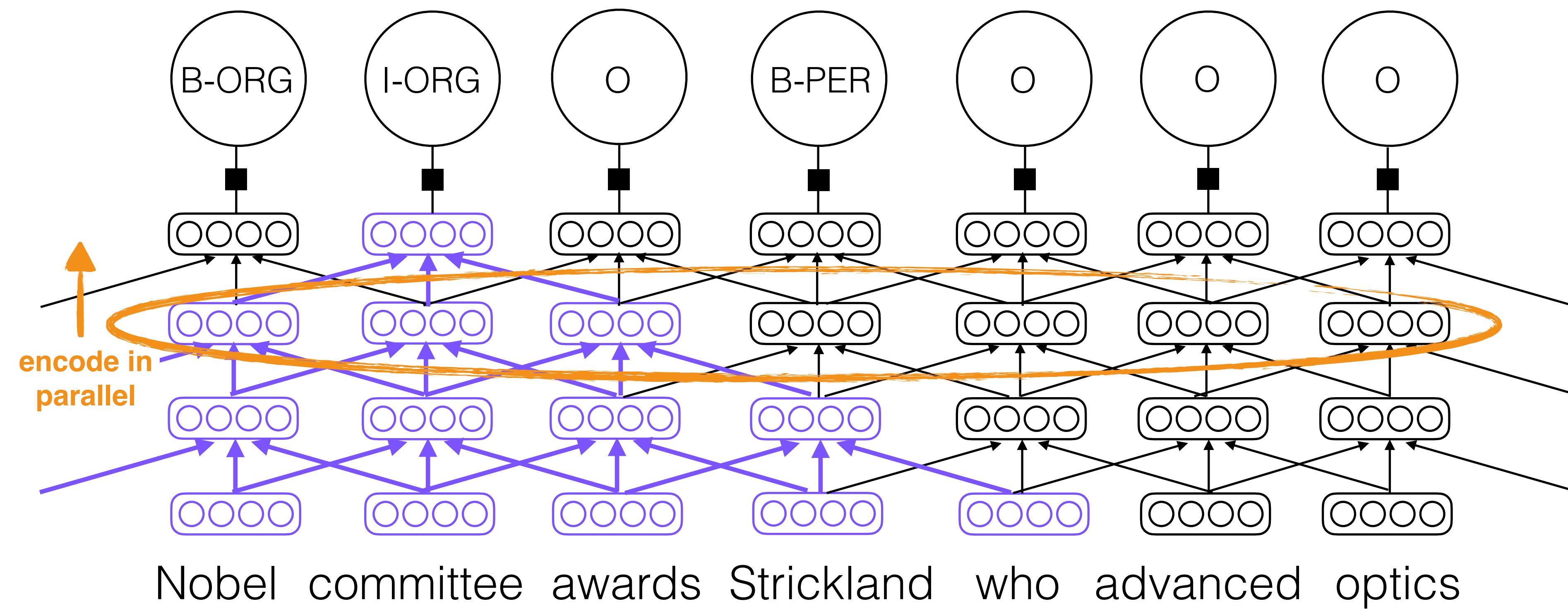
# Sequence labeling w/ CNNs



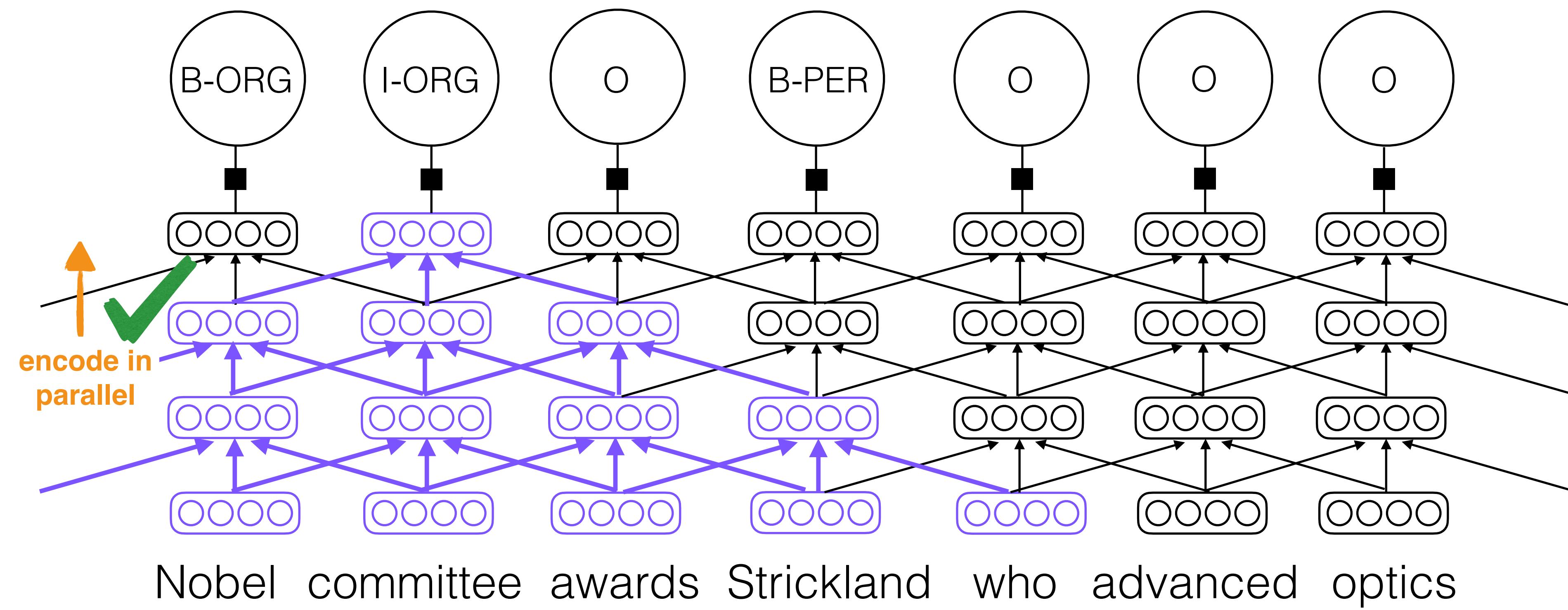
# Sequence labeling w/ CNNs



# Sequence labeling w/ CNNs

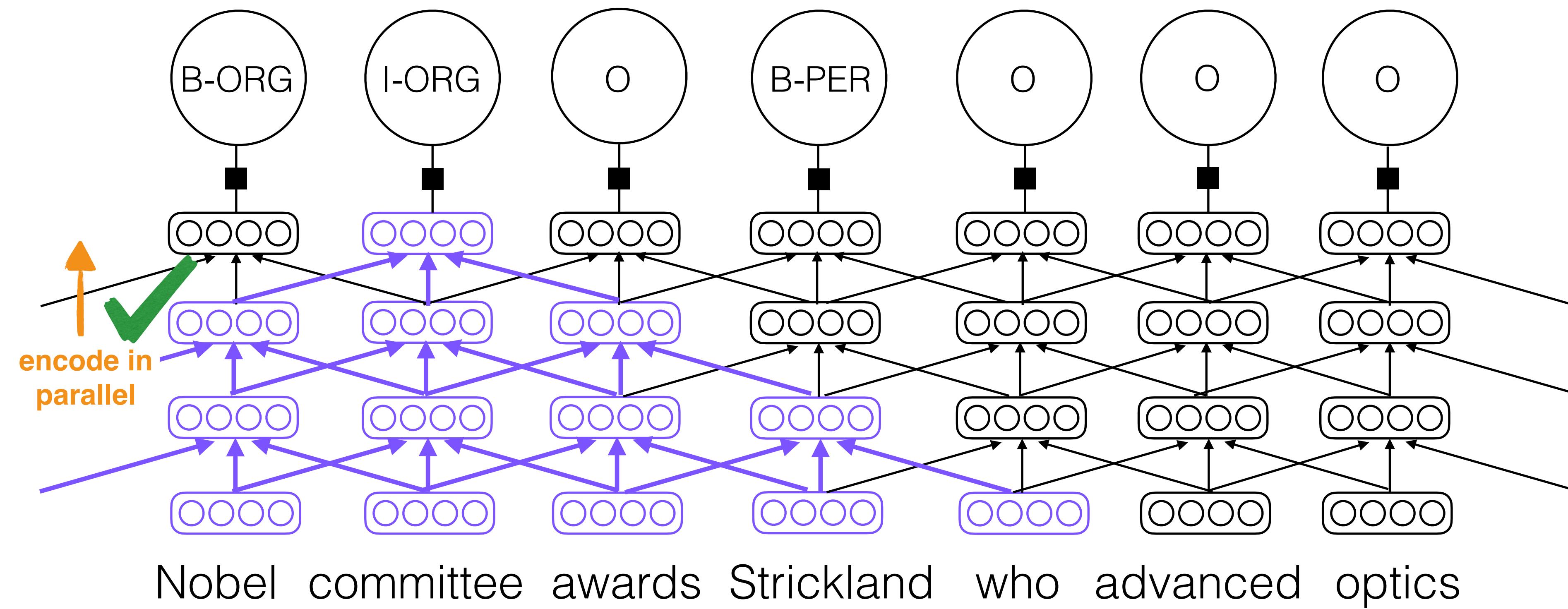


# Sequence labeling w/ CNNs



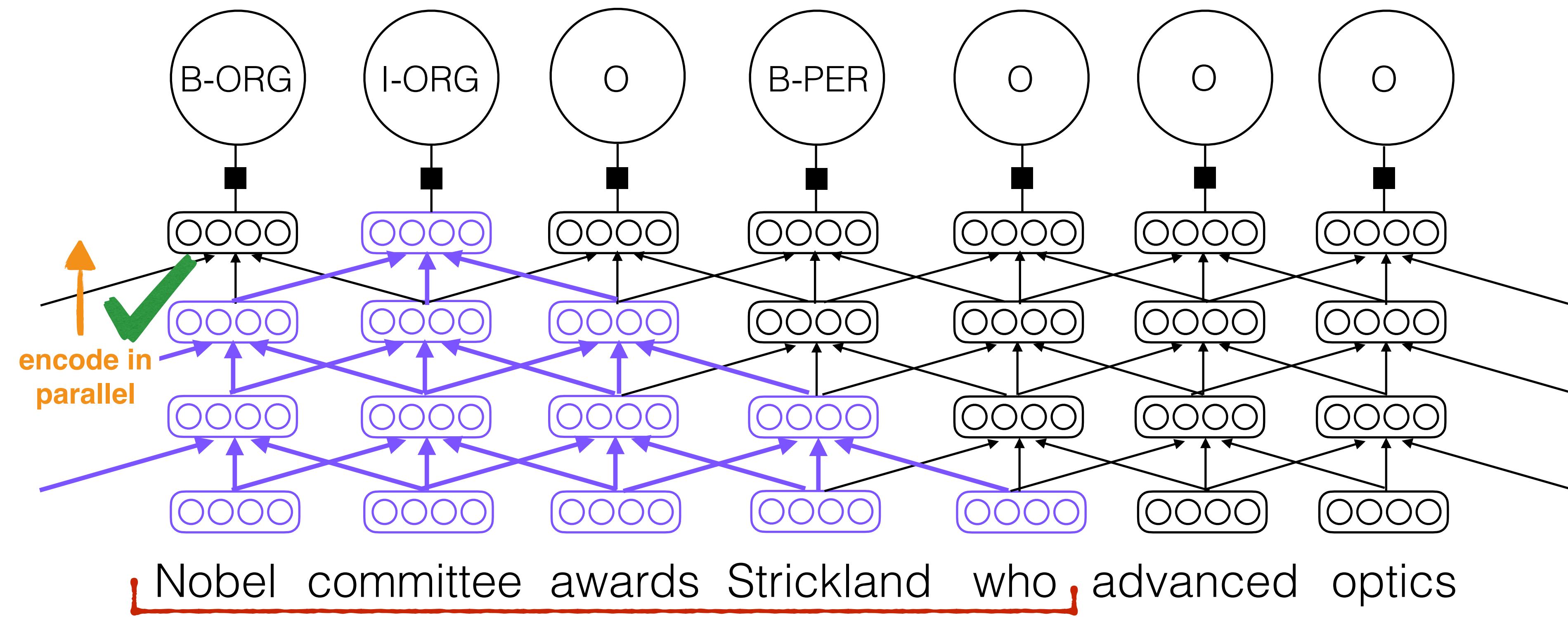
# Sequence labeling w/ CNNs

- Used for semantic role labeling, with poor results [[Collobert et al. 2011](#)].



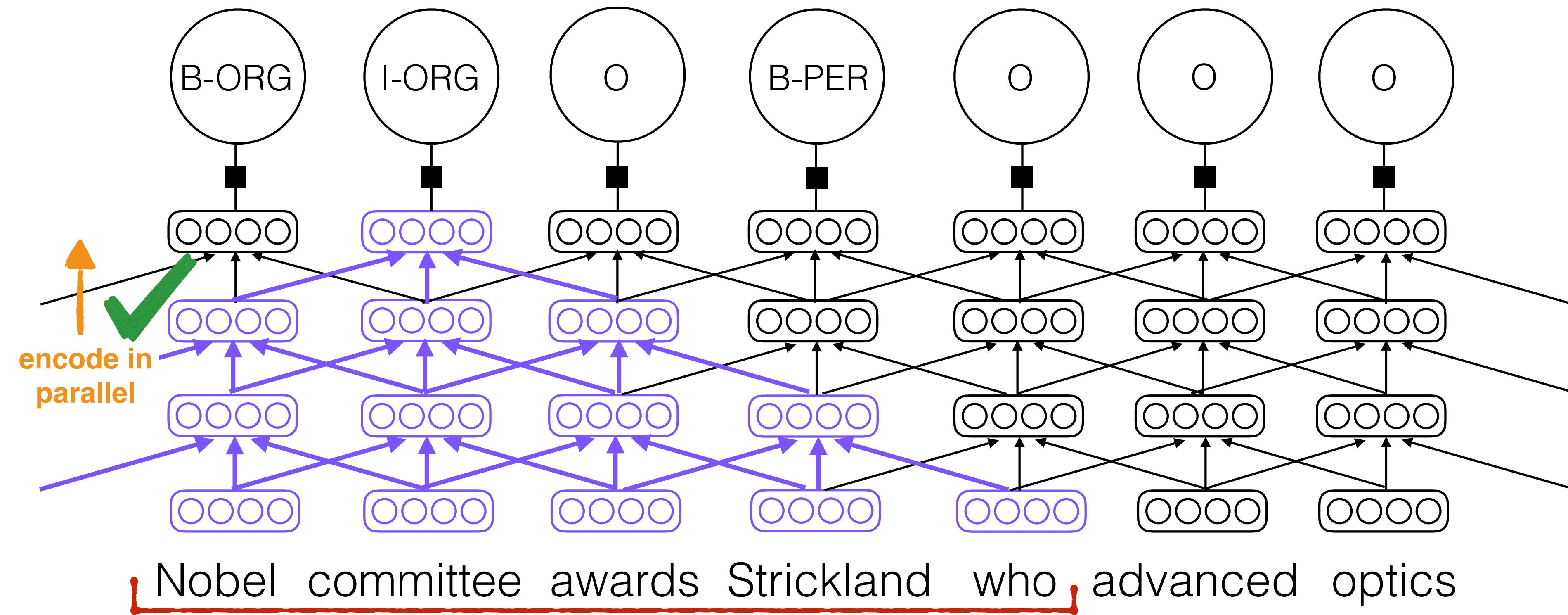
# Sequence labeling w/ CNNs

- Used for semantic role labeling, with poor results [[Collobert et al. 2011](#)].



# Sequence labeling w/ CNNs

- Used for semantic role labeling, with poor results [[Collobert et al. 2011](#)].
- Not enough context: amount of context grows **linearly** w/ number of layers.



# Sequence labeling w/ dilated CNNs



Nobel committee awards Strickland who advanced optics

# Sequence labeling w/ dilated CNNs

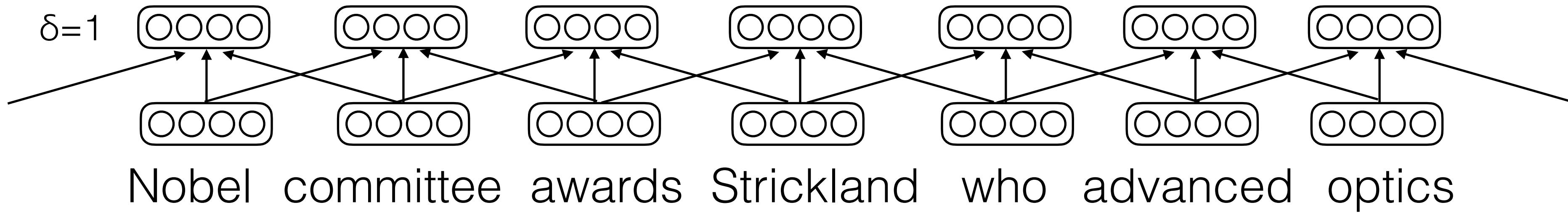
- Additional parameter: **dilation width  $\delta$**



Nobel committee awards Strickland who advanced optics

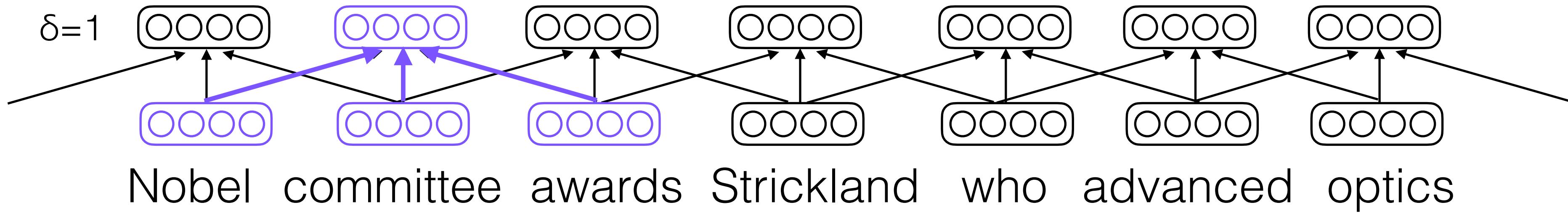
# Sequence labeling w/ dilated CNNs

- Additional parameter: **dilation width  $\delta$**



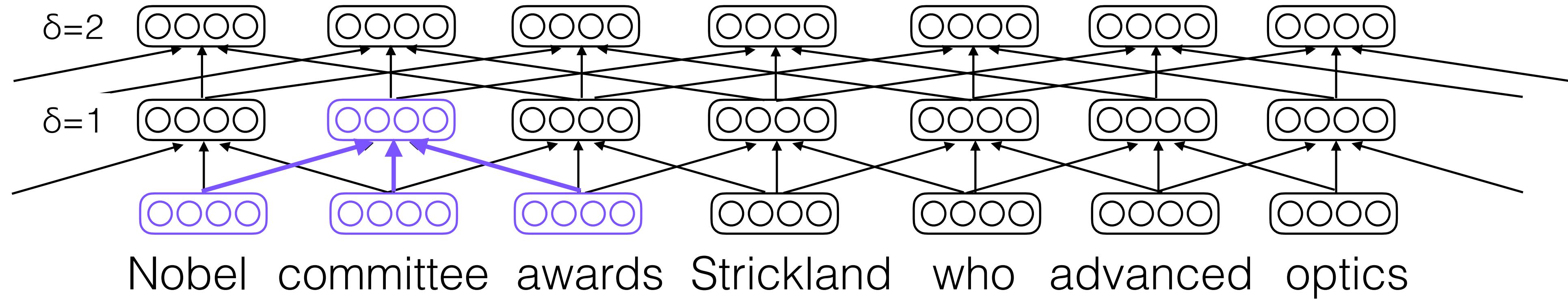
# Sequence labeling w/ dilated CNNs

- Additional parameter: **dilation width  $\delta$**



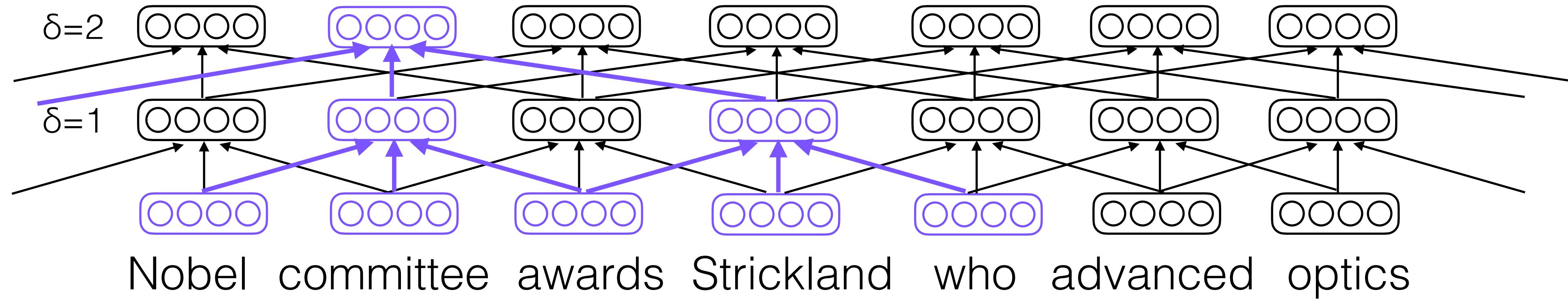
# Sequence labeling w/ dilated CNNs

- Additional parameter: **dilation width  $\delta$**



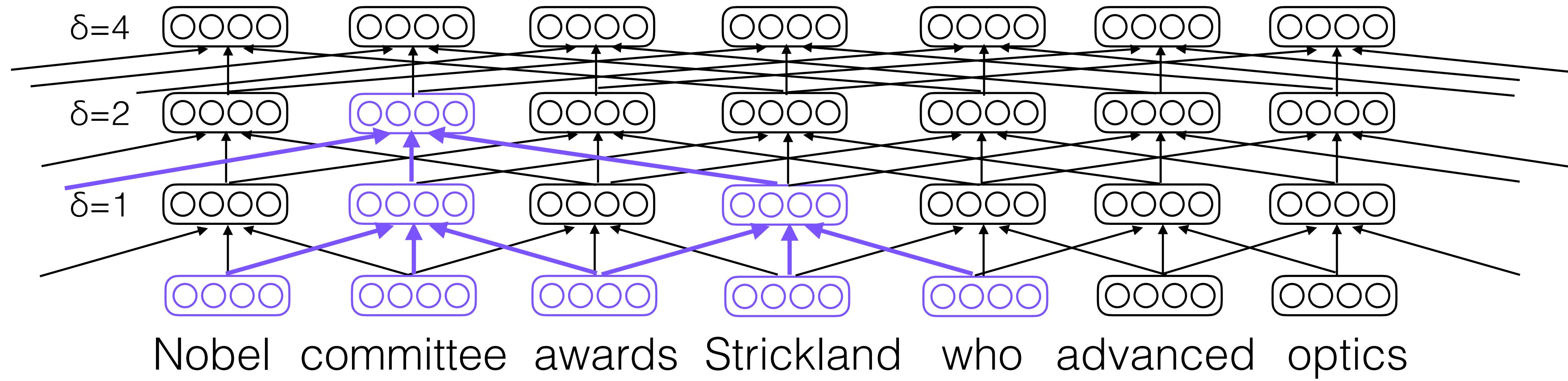
# Sequence labeling w/ dilated CNNs

- Additional parameter: **dilation width  $\delta$**



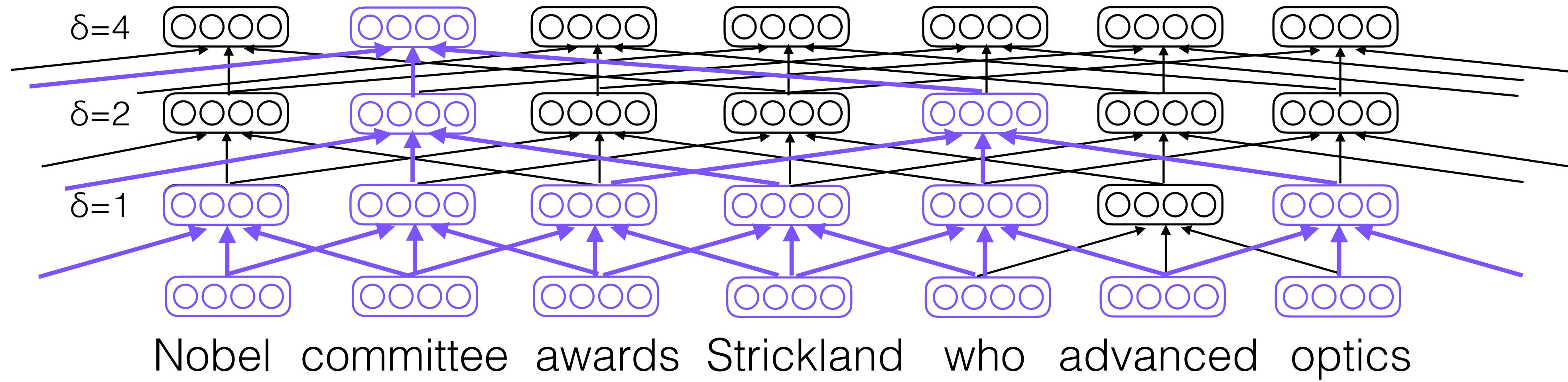
# Sequence labeling w/ dilated CNNs

- Additional parameter: **dilation width  $\delta$**



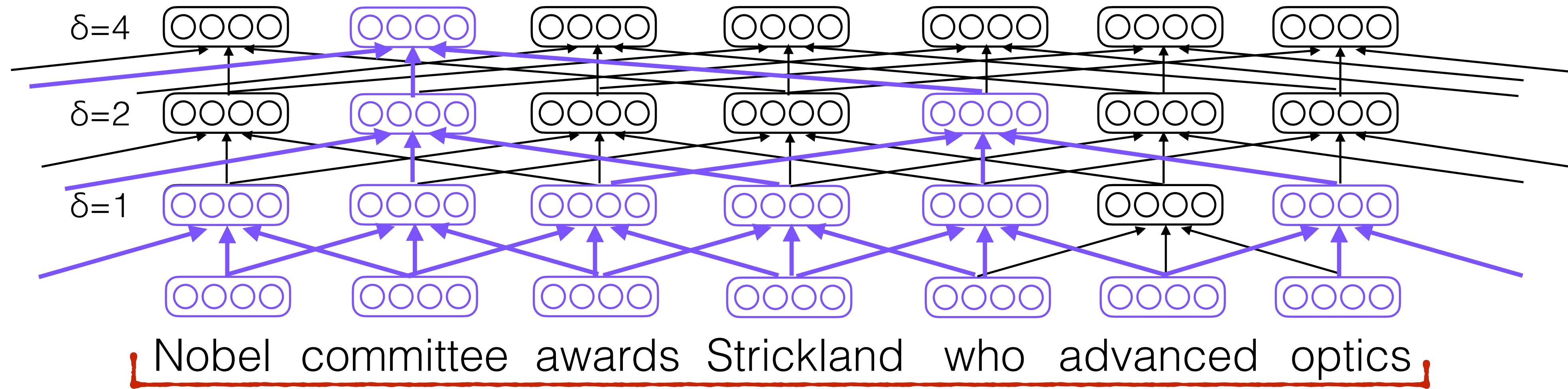
# Sequence labeling w/ dilated CNNs

- Additional parameter: **dilation width  $\delta$**



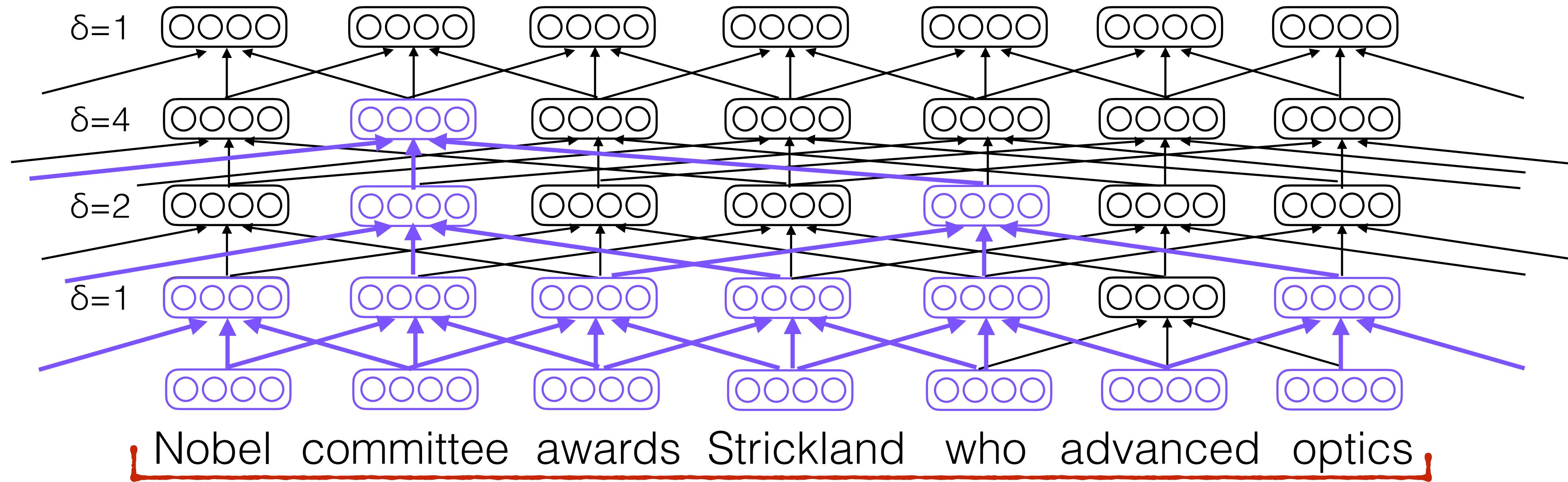
# Sequence labeling w/ dilated CNNs

- Additional parameter: **dilation width  $\delta$**



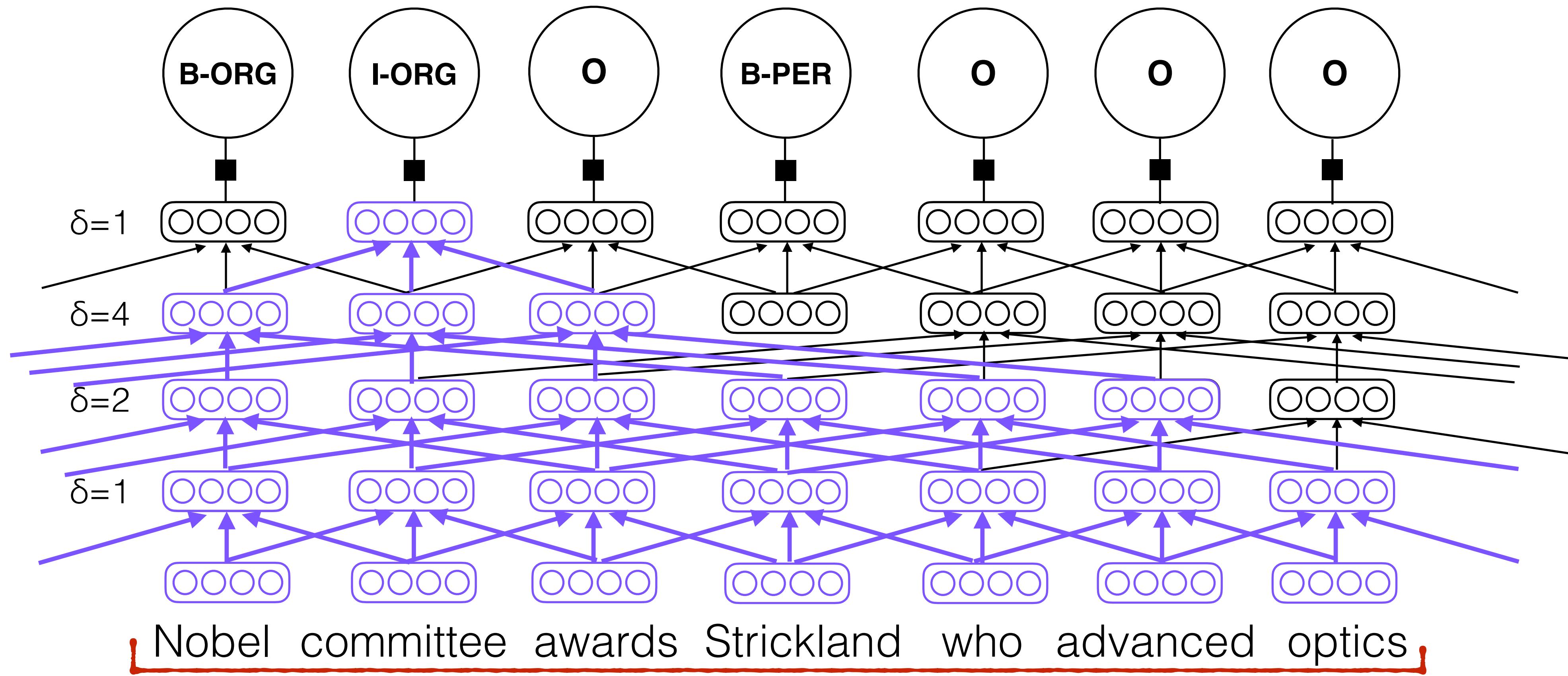
# Sequence labeling w/ dilated CNNs

- Additional parameter: **dilation width  $\delta$**



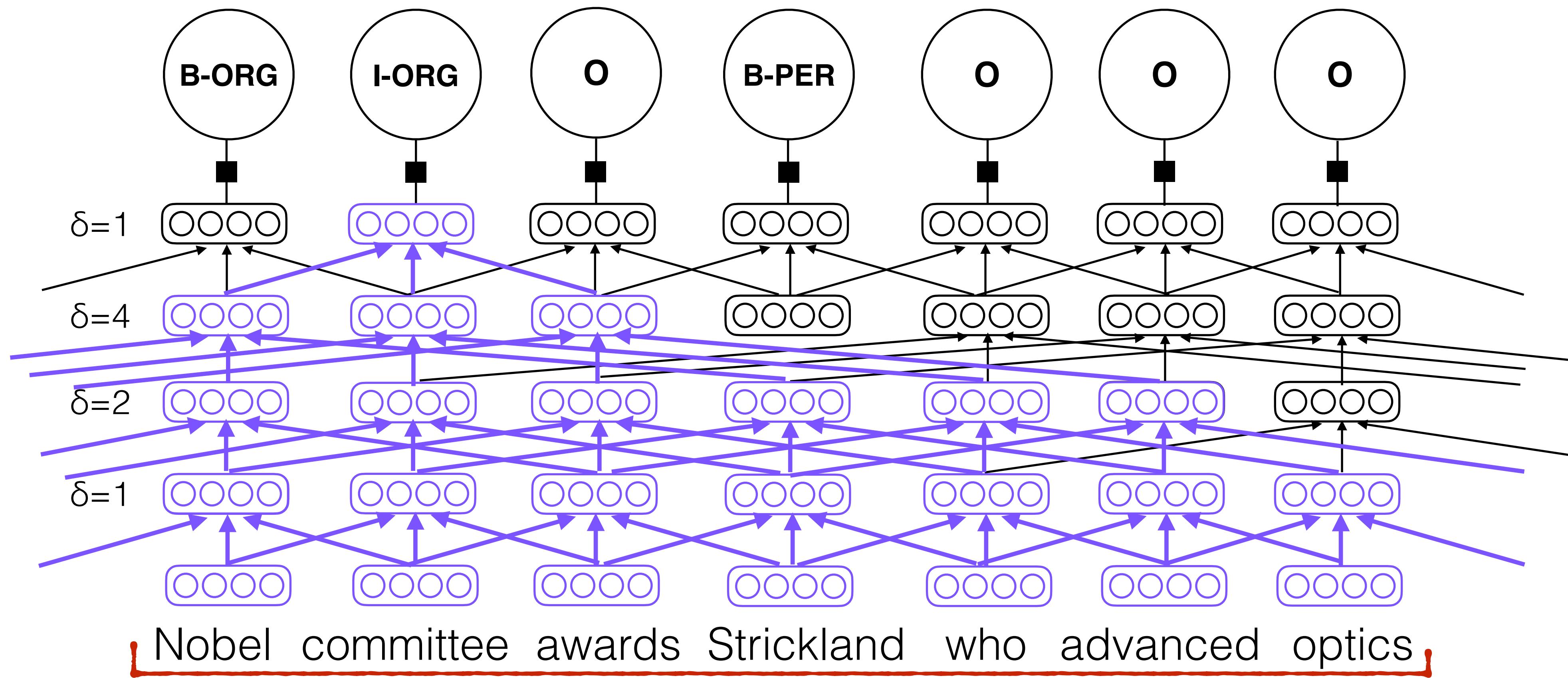
# Sequence labeling w/ dilated CNNs

- Additional parameter: **dilation width  $\delta$**



# Sequence labeling w/ dilated CNNs

- Additional parameter: **dilation width  $\delta$**
- Context window grows **exponentially** w/ number of layers.



# **Sequence labeling w/ dilated CNNs**

# Sequence labeling w/ dilated CNNs

- Why use a (dilated) CNN over a (bidirectional) LSTM?

# Sequence labeling w/ dilated CNNs

- Why use a (dilated) CNN over a (bidirectional) LSTM?
- Efficiency (on GPUs). Representations for every token in the sequence can be computed in parallel for CNN; linear dependence on sequence length for LSTM.

# Sequence labeling w/ dilated CNNs

- Why use a (dilated) CNN over a (bidirectional) LSTM?
- Efficiency (on GPUs). Representations for every token in the sequence can be computed in parallel for CNN; linear dependence on sequence length for LSTM.

	NER F1	
	sentence	document
Bi-LSTM-CRF	$90.4 \pm 0.1$	$90.6 \pm 0.2$
ID-CNN (ours)	$90.3 \pm 0.3$	$90.7 \pm 0.2$

# Sequence labeling w/ dilated CNNs

- Why use a (dilated) CNN over a (bidirectional) LSTM?
- Efficiency (on GPUs). Representations for every token in the sequence can be computed in parallel for CNN; linear dependence on sequence length for LSTM.

	NER F1	
	sentence	document
Bi-LSTM-CRF	$90.4 \pm 0.1$	$90.6 \pm 0.2$
ID-CNN (ours)	$90.3 \pm 0.3$	$90.7 \pm 0.2$

# Sequence labeling w/ dilated CNNs

- Why use a (dilated) CNN over a (bidirectional) LSTM?
- Efficiency (on GPUs). Representations for every token in the sequence can be computed in parallel for CNN; linear dependence on sequence length for LSTM.

	NER F1	
	sentence	document
Bi-LSTM-CRF	$90.4 \pm 0.1$	$90.6 \pm 0.2$
ID-CNN (ours)	$90.3 \pm 0.3$	$90.7 \pm 0.2$

*14x speed-up*

# Sequence labeling w/ dilated CNNs

- Why use a (dilated) CNN over a (bidirectional) LSTM?
- Efficiency (on GPUs). Representations for every token in the sequence can be computed in parallel for CNN; linear dependence on sequence length for LSTM.

	NER F1	
	sentence	document
Bi-LSTM-CRF	$90.4 \pm 0.1$	$90.6 \pm 0.2$
ID-CNN (ours)	$90.3 \pm 0.3$	$90.7 \pm 0.2$

*14x speed-up*

# Sequence labeling w/ dilated CNNs

- Why use a (dilated) CNN over a (bidirectional) LSTM?
- Efficiency (on GPUs). Representations for every token in the sequence can be computed in parallel for CNN; linear dependence on sequence length for LSTM.

	NER F1	
	sentence	document
Bi-LSTM-CRF	$90.4 \pm 0.1$	$90.6 \pm 0.2$
ID-CNN (ours)	$90.3 \pm 0.3$	$90.7 \pm 0.2$

*14x speed-up*      *8x speed-up*

# Character embeddings

# Character embeddings

- Character-level representations of words help to deal with UNKs.

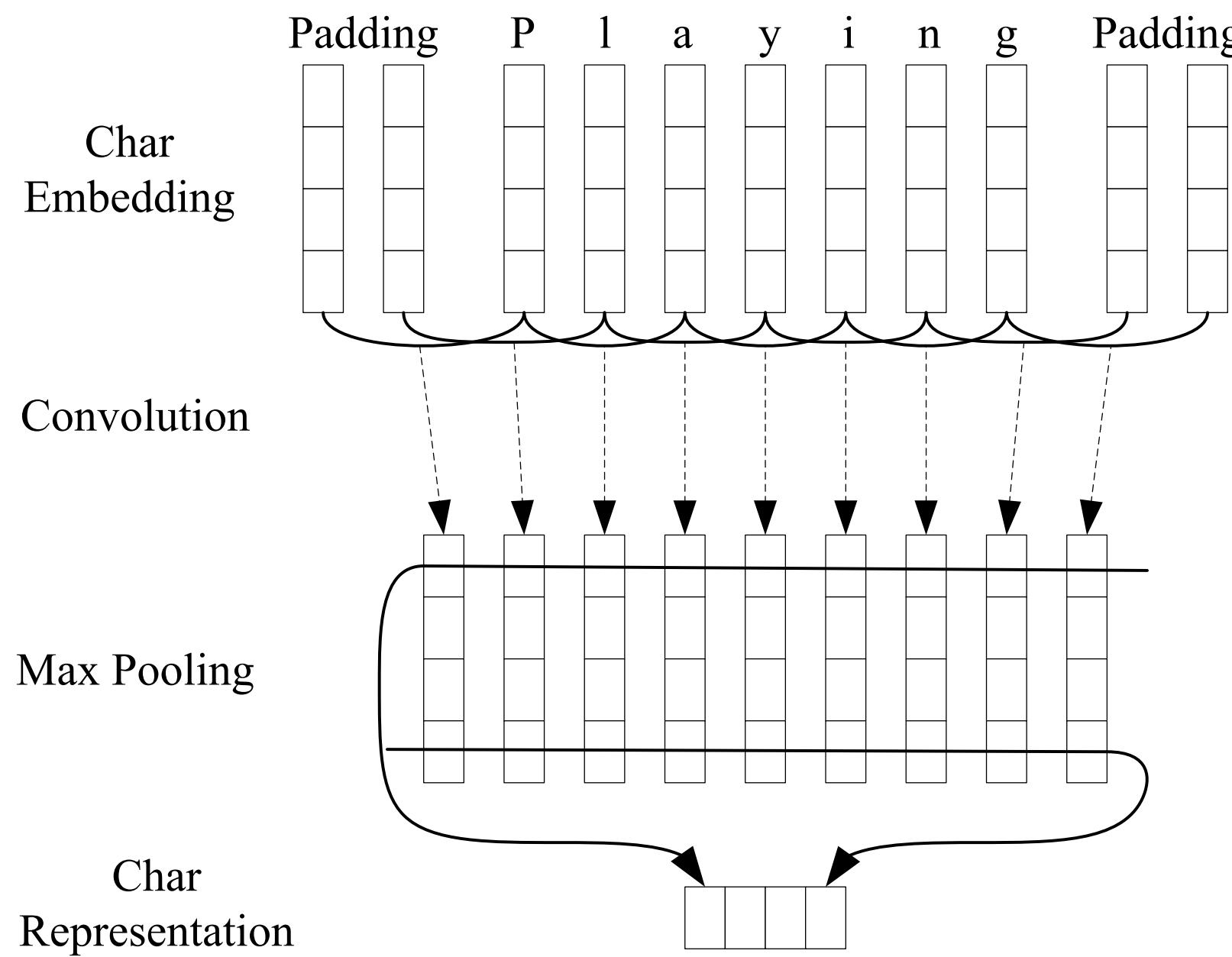
# Character embeddings

- Character-level representations of words help to deal with UNKs.
- Usually, CNNs + pooling are used to compose characters into word embeddings.

# Character embeddings

- Character-level representations of words help to deal with UNKs.
- Usually, CNNs + pooling are used to compose characters into word embeddings.

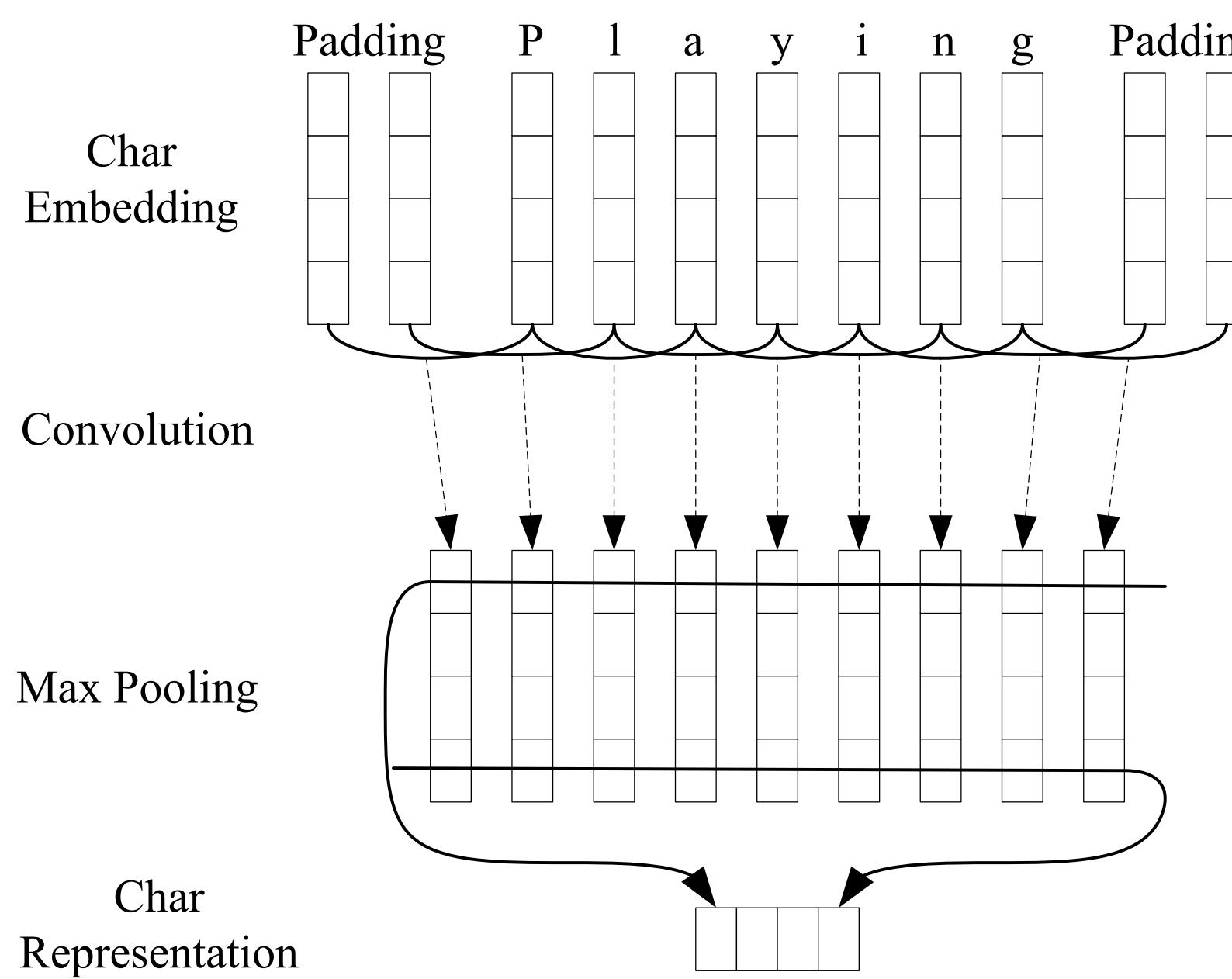
CNN + max pooling



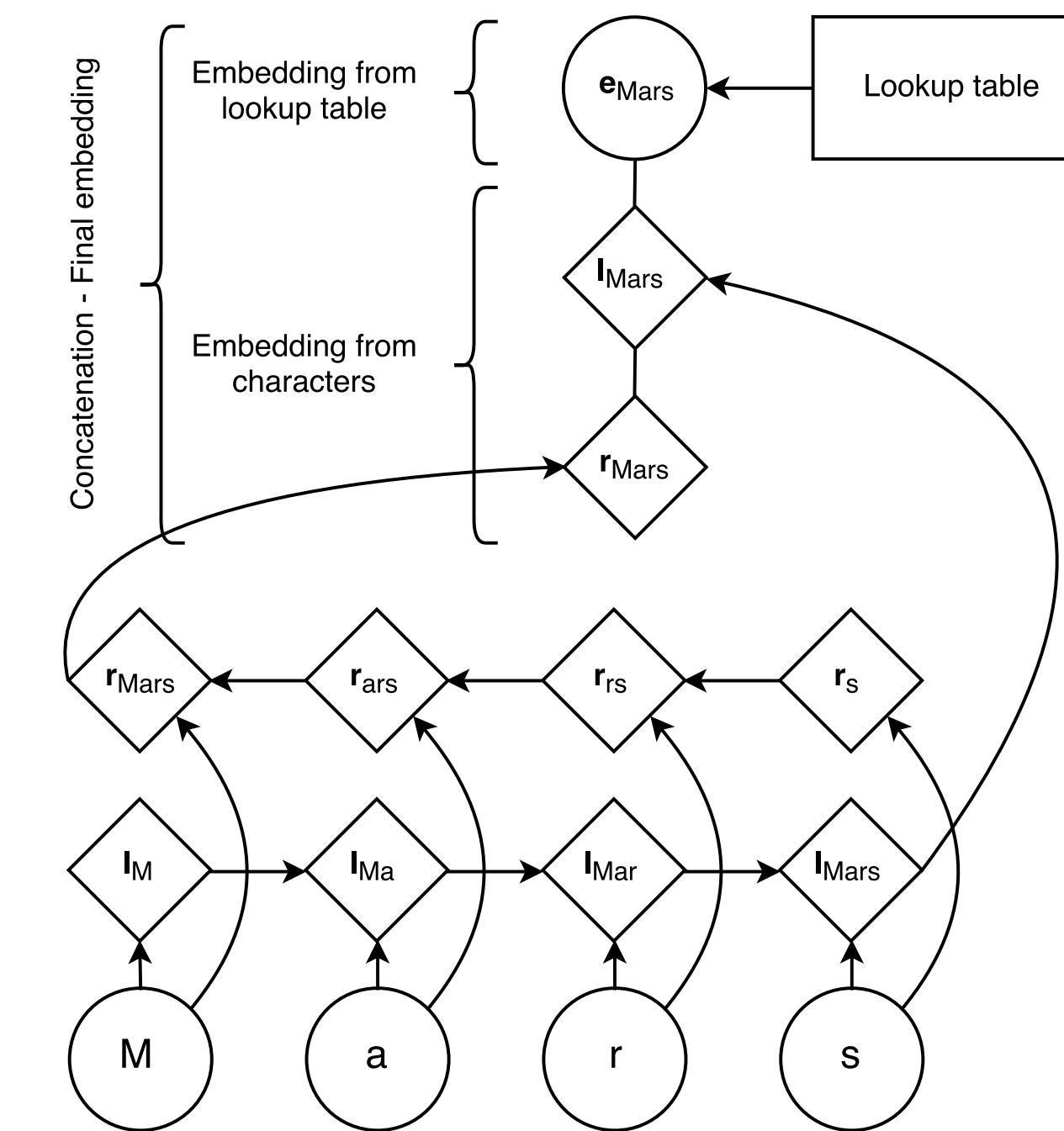
# Character embeddings

- Character-level representations of words help to deal with UNKs.
- Usually, CNNs + pooling are used to compose characters into word embeddings.

CNN + max pooling



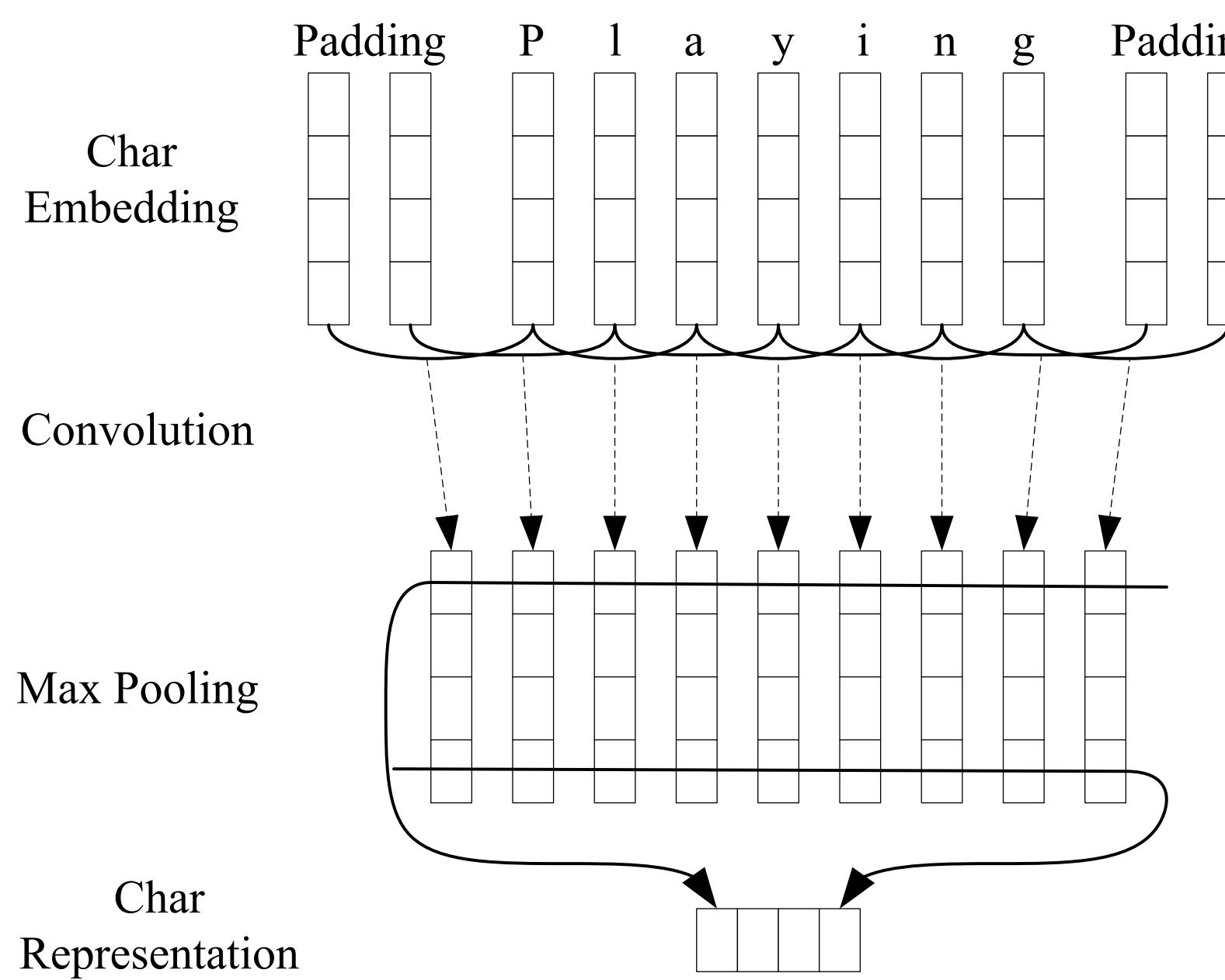
bidirectional LSTM



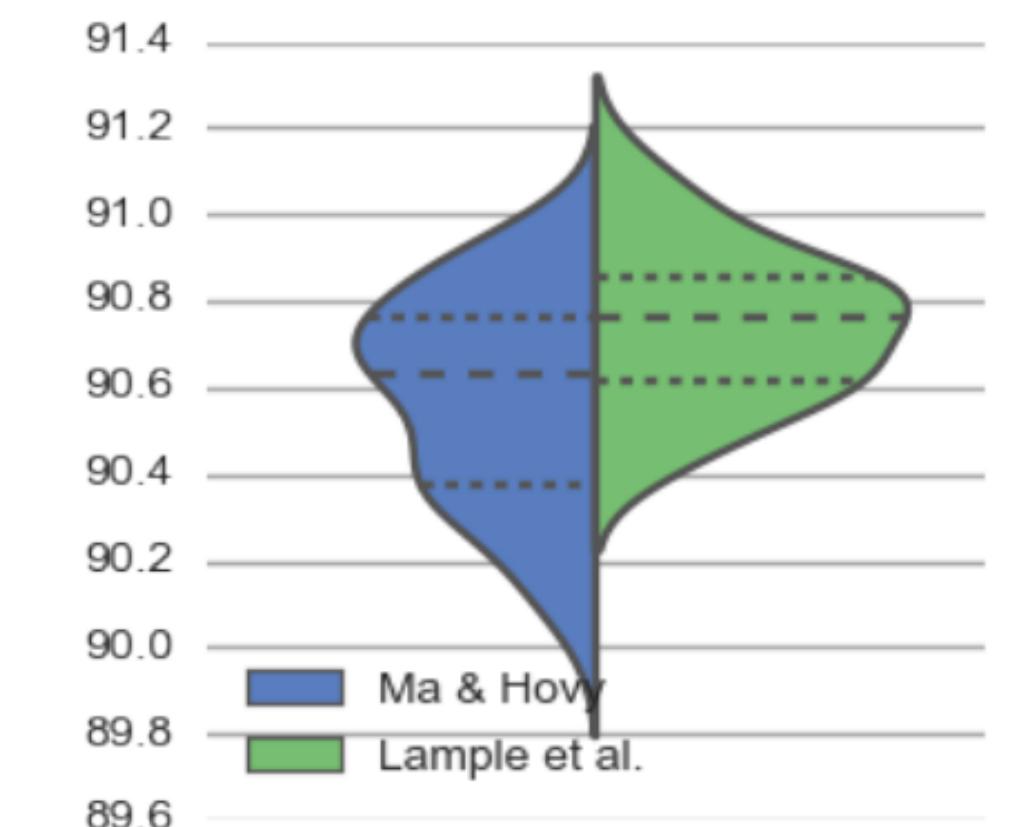
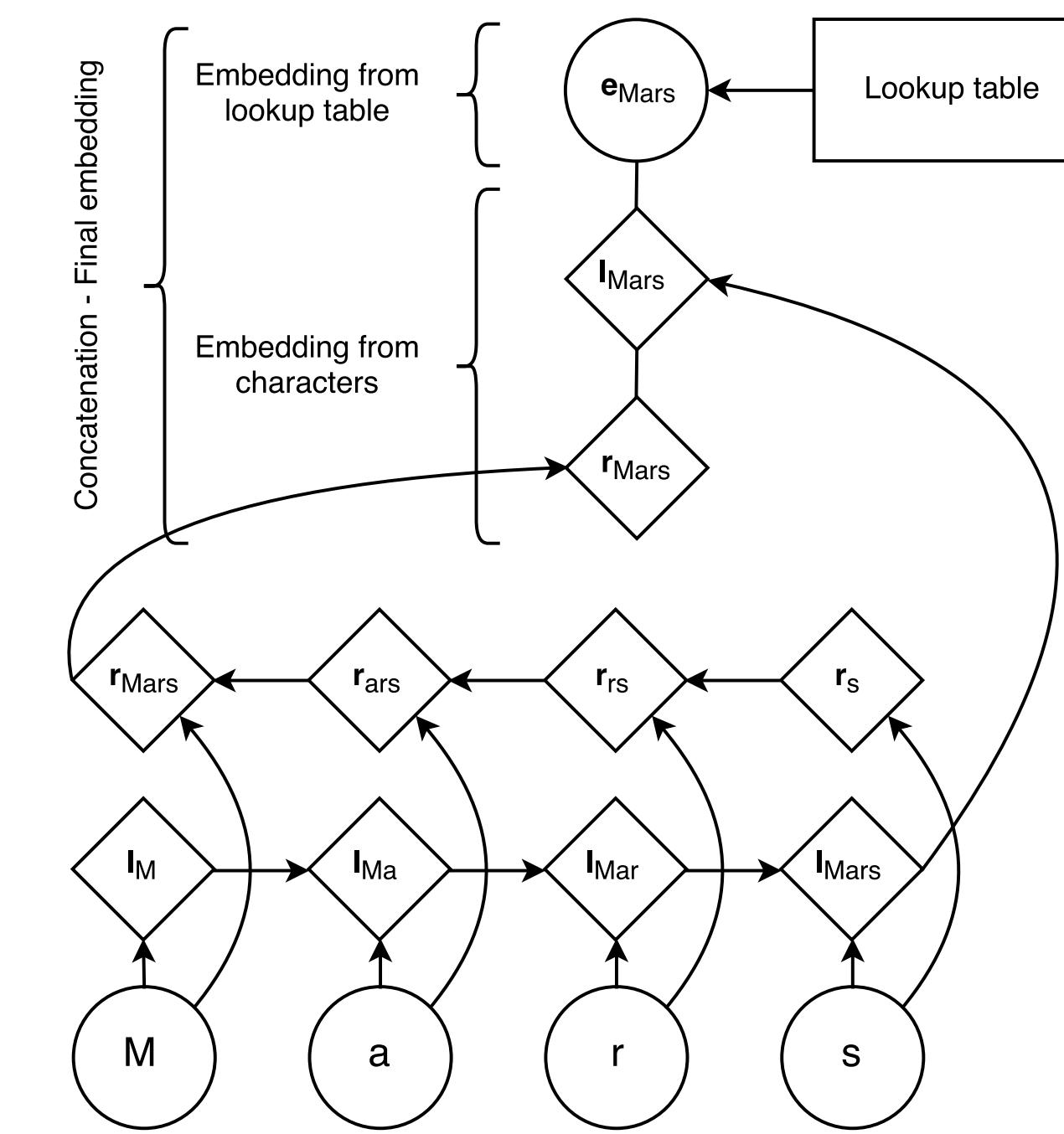
# Character embeddings

- Character-level representations of words help to deal with UNKs.
- Usually, CNNs + pooling are used to compose characters into word embeddings.

CNN + max pooling



bidirectional LSTM



# Multilingual part-of-speech tagging

# Multilingual part-of-speech tagging

- Many UNKs in morphologically-rich languages like Czech, Hungarian, Turkish

# Multilingual part-of-speech tagging

- Many UNKs in morphologically-rich languages like Czech, Hungarian, Turkish
  - 250,000 word corpus of Hungarian has > 2x as many types as a similarly sized corpus of English

# Multilingual part-of-speech tagging

- Many UNKs in morphologically-rich languages like Czech, Hungarian, Turkish
  - 250,000 word corpus of Hungarian has > 2x as many types as a similarly sized corpus of English
  - 10 million word corpus of Turkish contains 4x as many types

# Multilingual part-of-speech tagging

- Many UNKs in morphologically-rich languages like Czech, Hungarian, Turkish
  - 250,000 word corpus of Hungarian has > 2x as many types as a similarly sized corpus of English
  - 10 million word corpus of Turkish contains 4x as many types
- Information coded in morphology

# Multilingual part-of-speech tagging

- Many UNKs in morphologically-rich languages like Czech, Hungarian, Turkish
  - 250,000 word corpus of Hungarian has > 2x as many types as a similarly sized corpus of English
  - 10 million word corpus of Turkish contains 4x as many types
- Information coded in morphology

partilerindeydi

partisindeydiler

# Multilingual part-of-speech tagging

- Many UNKs in morphologically-rich languages like Czech, Hungarian, Turkish
  - 250,000 word corpus of Hungarian has > 2x as many types as a similarly sized corpus of English
  - 10 million word corpus of Turkish contains 4x as many types
- Information coded in morphology

partilerindeydi

partisindeydiler

*he/she/they<sub>(sing)</sub> were/was at their<sub>(plur)</sub> party*

# Multilingual part-of-speech tagging

- Many UNKs in morphologically-rich languages like Czech, Hungarian, Turkish
  - 250,000 word corpus of Hungarian has > 2x as many types as a similarly sized corpus of English
  - 10 million word corpus of Turkish contains 4x as many types
- Information coded in morphology

partilerindeydi

*he/she/they<sub>(sing)</sub> were/was at their<sub>(plur)</sub> party*

partisindeydiler

*they<sub>(plur)</sub> were at his/her/their<sub>(sing)</sub> party*

# Multilingual part-of-speech tagging

- Many UNKs in morphologically-rich languages like Czech, Hungarian, Turkish
  - 250,000 word corpus of Hungarian has > 2x as many types as a similarly sized corpus of English
  - 10 million word corpus of Turkish contains 4x as many types
- Information coded in morphology

partilerindeydi

party

*he/she/they<sub>(sing)</sub> were/was at their<sub>(plur)</sub> party*

partisindeydiler

party

*they<sub>(plur)</sub> were at his/her/their<sub>(sing)</sub> party*

# Multilingual part-of-speech tagging

- Many UNKs in morphologically-rich languages like Czech, Hungarian, Turkish
  - 250,000 word corpus of Hungarian has > 2x as many types as a similarly sized corpus of English
  - 10 million word corpus of Turkish contains 4x as many types
- Information coded in morphology

in  
partilerindeydi  
party

he/she/they<sub>(sing)</sub> were/was at their<sub>(plur)</sub> party

in  
partisindeydiler  
party

they<sub>(plur)</sub> were at his/her/their<sub>(sing)</sub> party

# Multilingual part-of-speech tagging

- Many UNKs in morphologically-rich languages like Czech, Hungarian, Turkish
  - 250,000 word corpus of Hungarian has > 2x as many types as a similarly sized corpus of English
  - 10 million word corpus of Turkish contains 4x as many types
- Information coded in morphology

in she/he/they<sub>(sing)</sub> was  
partilerindeydi  
party their  
he/she/they<sub>(sing)</sub> were/was at their<sub>(plur)</sub> party

in they<sub>(plur)</sub> were  
partisindeydiler  
party his/her/their<sub>(sing)</sub>  
they<sub>(plur)</sub> were at his/her/their<sub>(sing)</sub> party

# Multilingual part-of-speech tagging

- Many UNKs in morphologically-rich languages like Czech, Hungarian, Turkish
  - 250,000 word corpus of Hungarian has > 2x as many types as a similarly sized corpus of English
  - 10 million word corpus of Turkish contains 4x as many types
- Information coded in morphology

in she/he/they<sub>(sing)</sub> was  
partilerindeydi

party their

he/she/they<sub>(sing)</sub> were/was at their<sub>(plur)</sub> party

in they<sub>(plur)</sub> were  
partisindeydiler

party his/her/their<sub>(sing)</sub>

they<sub>(plur)</sub> were at his/her/their<sub>(sing)</sub> party

Yerdeki **izin** temizlenmesi gerek.  
The trace on the floor should be cleaned.

iz + Noun+A3sg+Pnon+Gen

# Multilingual part-of-speech tagging

# Multilingual part-of-speech tagging

- In non-word-space languages like Chinese, word segmentation is either applied before tagging or performed jointly.

# Multilingual part-of-speech tagging

- In non-word-space languages like Chinese, word segmentation is either applied before tagging or performed jointly.

姚明进入总决赛

*Yao Ming reaches the finals*

# Multilingual part-of-speech tagging

- In non-word-space languages like Chinese, word segmentation is either applied before tagging or performed jointly.

姚明进入总决赛

Yao Ming reaches the finals

姚明 进入 总决赛

YaoMing reaches finals

# Multilingual part-of-speech tagging

- In non-word-space languages like Chinese, word segmentation is either applied before tagging or performed jointly.

姚明进入总决赛

Yao Ming reaches the finals

CTB

姚明 进入 总决赛

YaoMing reaches finals

# Multilingual part-of-speech tagging

- In non-word-space languages like Chinese, word segmentation is either applied before tagging or performed jointly.

姚明进入总决赛

Yao Ming reaches the finals

CTB

姚明 进入 总决赛

YaoMing reaches finals

姚 明 进入 总 决赛

Yao Ming reaches overall finals

# Multilingual part-of-speech tagging

- In non-word-space languages like Chinese, word segmentation is either applied before tagging or performed jointly.

姚明进入总决赛

Yao Ming reaches the finals

CTB

姚明 进入 总决赛

YaoMing reaches finals

Peking U.

姚 明 进 入 总 决 赛

Yao Ming reaches overall finals

# Multilingual part-of-speech tagging

- In non-word-space languages like Chinese, word segmentation is either applied before tagging or performed jointly.
- UNKs are difficult: majority of unknown words are common nouns and verbs due to compounding

姚明进入总决赛

Yao Ming reaches the finals

CTB

姚明 进入 总决赛

YaoMing reaches finals

Peking U.

姚 明 进 入 总 决 赛

Yao Ming reaches overall finals

# Multilingual part-of-speech tagging

- In non-word-space languages like Chinese, word segmentation is either applied before tagging or performed jointly.
- UNKs are difficult: majority of unknown words are common nouns and verbs due to compounding

姚明进入总决赛

Yao Ming reaches the finals

CTB

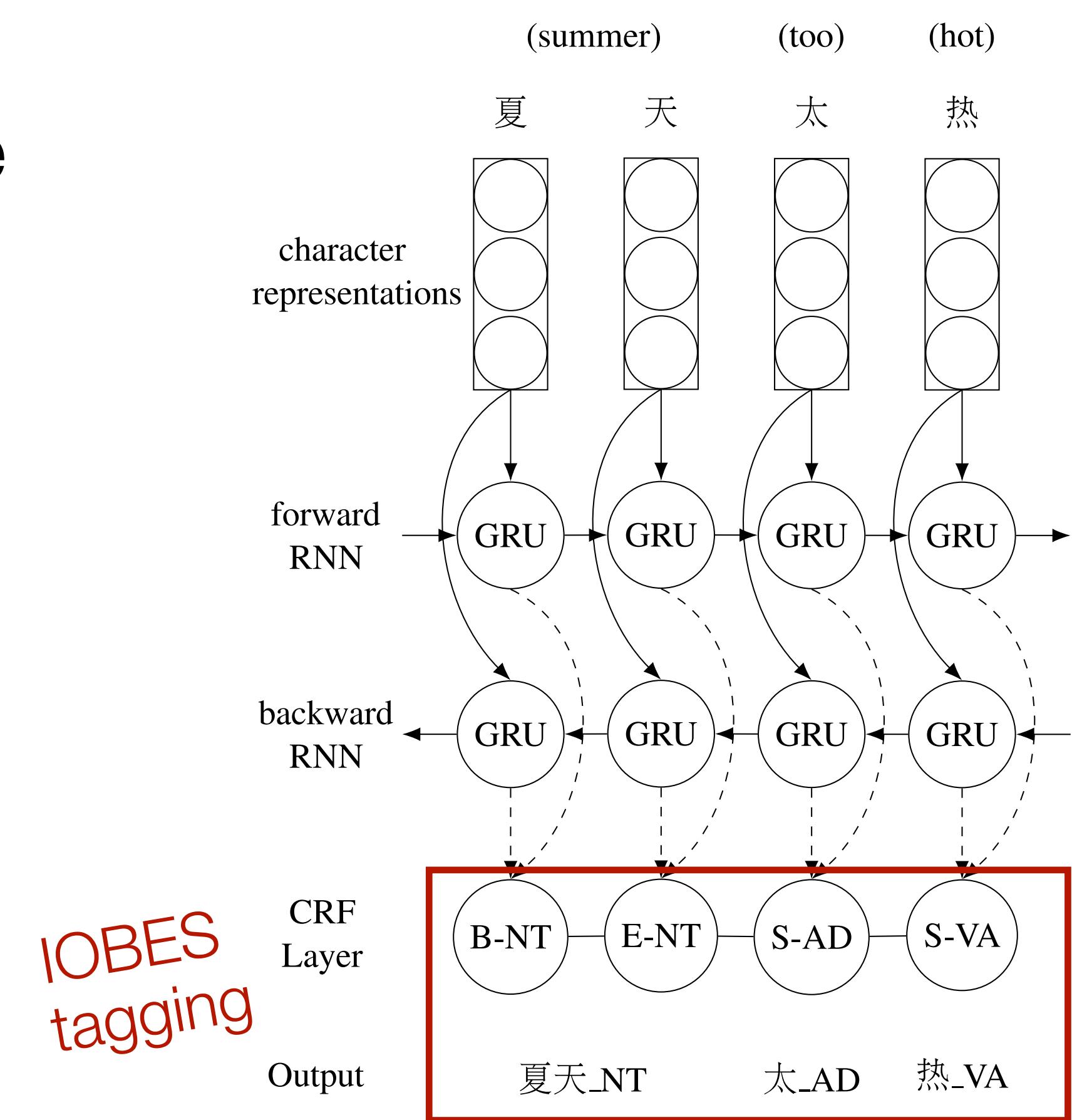
姚明 进入 总决赛

YaoMing reaches finals

Peking U.

姚 明 进入 总 决赛

Yao Ming reaches overall finals



# Multilingual part-of-speech tagging

# Multilingual part-of-speech tagging

- Universal POS tags [[Petrov et al. 2012](#)] provide a cross-lingual tag set.

Language	Source	# Tags
Arabic	PADT/CoNLL07 (Hajič et al., 2004)	21
Basque	Basque3LB/CoNLL07 (Aduriz et al., 2003)	64
Bulgarian	BTB/CoNLL06 (Simov et al., 2002)	54
Catalan	CESS-ECE/CoNLL07 (Martí et al., 2007)	54
Chinese	Penn ChineseTreebank 6.0 (Palmer et al., 2007)	34
Chinese	Sinica/CoNLL07 (Chen et al., 2003)	294
Czech	PDT/CoNLL07 (Böhmová et al., 2003)	63
Danish	DDT/CoNLL06 (Kromann et al., 2003)	25
Dutch	Alpino/CoNLL06 (Van der Beek et al., 2002)	12
English	PennTreebank (Marcus et al., 1993)	45
French	FrenchTreebank (Abeillé et al., 2003)	30
German	Tiger/CoNLL06 (Brants et al., 2002)	54
German	Negra (Skut et al., 1997)	54
Greek	GDT/CoNLL07 (Prokopidis et al., 2005)	38
Hungarian	Szeged/CoNLL07 (Csendes et al., 2005)	43
Italian	ISST/CoNLL07 (Montemagni et al., 2003)	28
Japanese	Verbmobil/CoNLL06 (Kawata and Bartels, 2000)	80
Japanese	Kyoto4.0 (Kurohashi and Nagao, 1997)	42
Korean	Sejong ( <a href="http://www.sejong.or.kr">http://www.sejong.or.kr</a> )	187
Portuguese	Floresta Sintá(c)tica/CoNLL06 (Afonso et al., 2002)	22
Russian	SynTagRus-RNC (Boguslavsky et al., 2002)	11
Slovene	SDT/CoNLL06 (Džeroski et al., 2006)	29
Spanish	Ancora-Cast3LB/CoNLL06 (Civit and Martí, 2004)	47
Swedish	Talbanken05/CoNLL06 (Nivre et al., 2006)	41
Turkish	METU-Sabancı/CoNLL07 (Oflazer et al., 2003)	31

# Multilingual part-of-speech tagging

- **Universal POS tags** [[Petrov et al. 2012](#)] provide a cross-lingual tag set.
  - Coarse grained: 16 tags

<b>open class</b>	<b>closed class</b>	<b>other</b>
ADJ	ADP	PUNCT
ADV	AUX	SYM
INTJ	CCONJ	X
NOUN	DET	
PROPN	NUM	
VERB	PART	
	PRON	
	SCONJ	

# Multilingual part-of-speech tagging

- **Universal POS tags** [[Petrov et al. 2012](#)] provide a cross-lingual tag set.
  - Coarse grained: 16 tags
  - Finer-grained analysis split off into morphological tags (case, gender, ...)

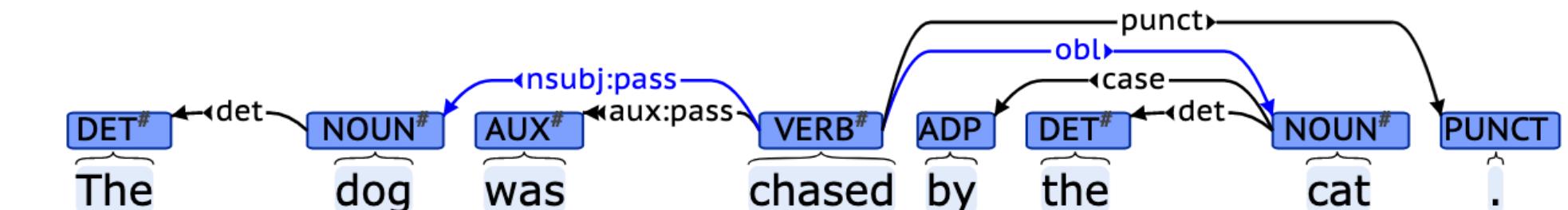
<b>open class</b>	<b>closed class</b>	<b>other</b>
ADJ	ADP	PUNCT
ADV	AUX	SYM
INTJ	CCONJ	X
NOUN	DET	
PROPN	NUM	
VERB	PART	
	PRON	
	SCONJ	

# Multilingual part-of-speech tagging

- Universal POS tags [Petrov et al. 2012] provide a cross-lingual tag set.
  - Coarse grained: 16 tags
  - Finer-grained analysis split off into morphological tags (case, gender, ...)

open class	closed class	other
ADJ	ADP	PUNCT
ADV	AUX	SYM
INTJ	CCONJ	X
NOUN	DET	
PROPN	NUM	
VERB	PART	
PRON		
SCONJ		

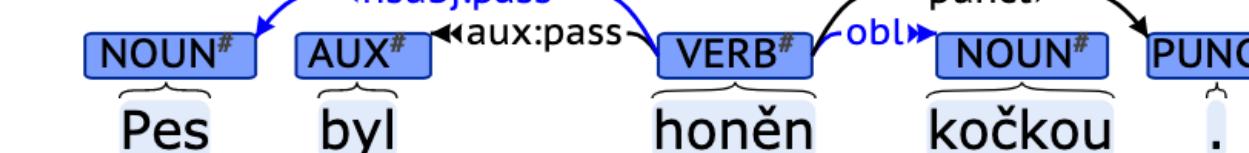
English:



Bulgarian:



Czech:



Swedish:



# Announcements

- **Project 2 released today after class:** sequence labeling.
  - Due: October 16.
  - You will implement part-of-speech taggers for English and Norwegian:
    - HMM, BiLSTM, and BiLSTM-CRF.
- Friday's recitation will be an overview of P2.