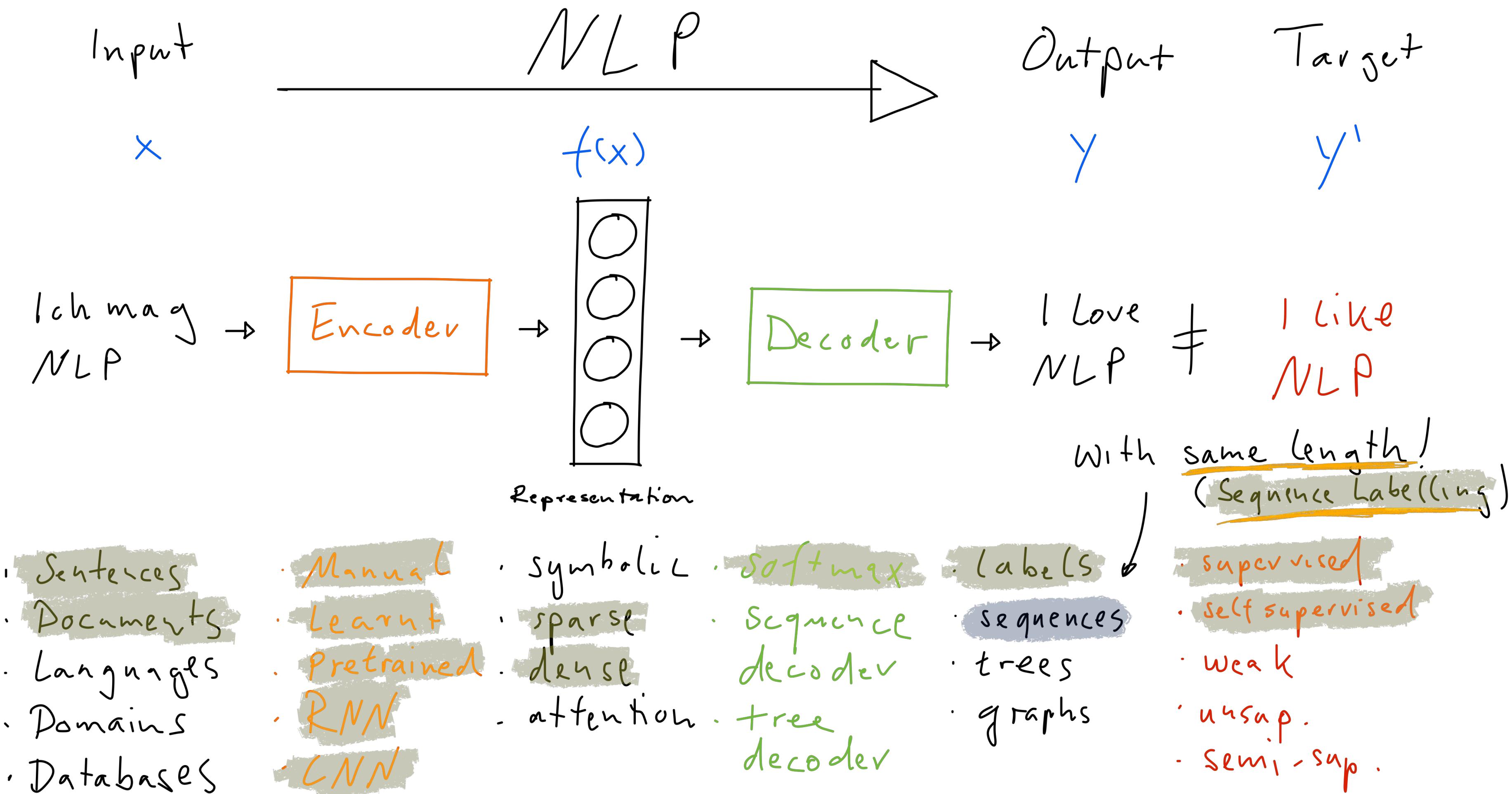


# Sequence Labelling

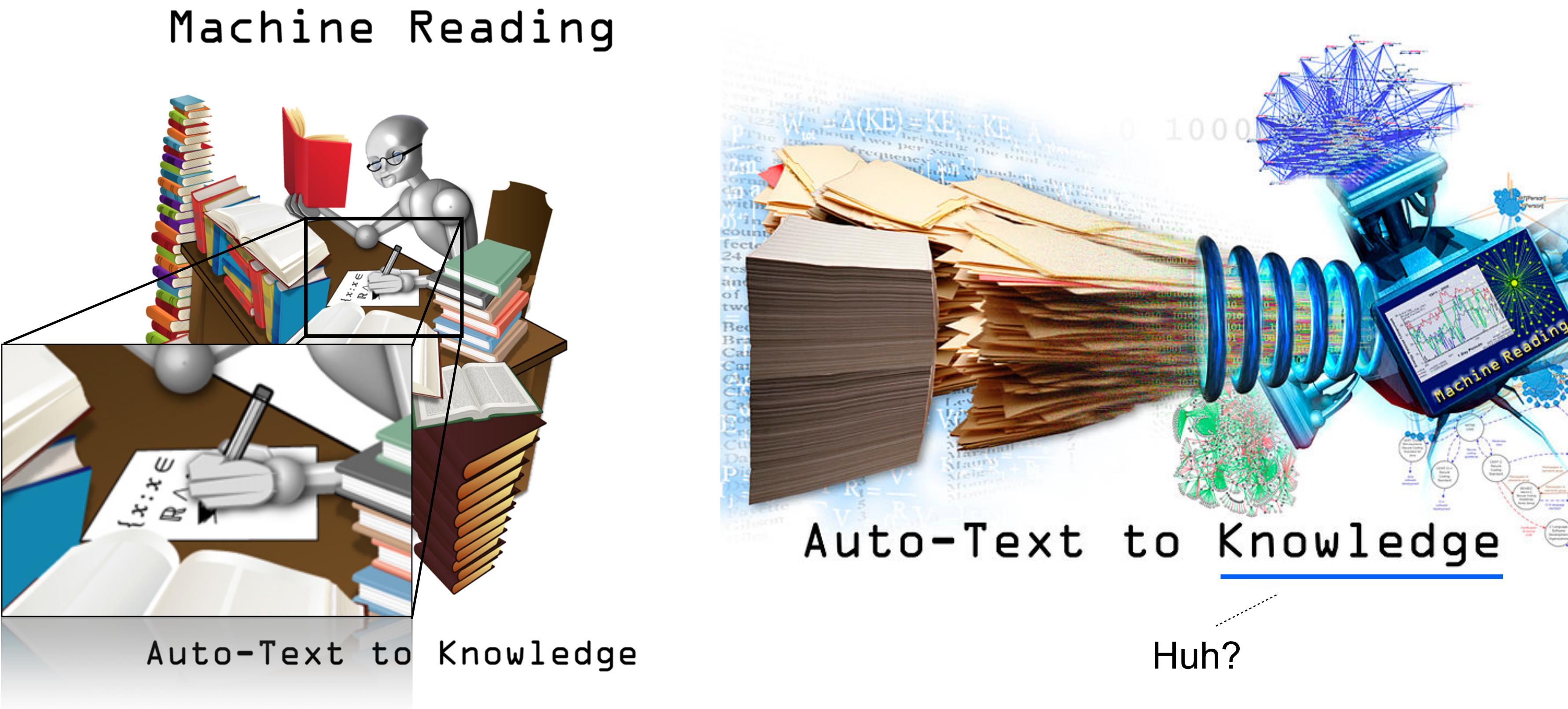
Tim Rocktäschel & **Sebastian Riedel**  
COMP0087 Natural Language Processing



# NLP in a Nutshell



# Machine Reading



# Machine Reading

andrew mccallum

Web Images Maps Shopping More Search tools

About 4,380,000 results (0.20 seconds)

Cookies help us deliver our services. By using our services, you agree to our use of cookies.  
OK Learn more

[Andrew McCallum Homepage](#)  
[www.cs.umass.edu/~mccallum](http://www.cs.umass.edu/~mccallum) ▾  
Machine learning, text and information retrieval and extraction, reinforcement learning.  
[Andrew McCallum Publications](#) - [Andrew McCallum Bio](#) - [People](#) - [Teaching](#)

[Andrew McCallum - London Metropolitan University](#)  
[www.londonmet.ac.uk/faculties/faculty-of...k.../andrew-mccallum/](http://www.londonmet.ac.uk/faculties/faculty-of...k.../andrew-mccallum/) ▾  
Andrew taught English in London secondary schools for 15 years before coming to London Met in 2008. He is course tutor for the PGCE in Secondary English ...

[Andrew McCallum - Wikipedia, the free encyclopedia](#)  
[en.wikipedia.org/wiki/Andrew\\_McCallum](http://en.wikipedia.org/wiki/Andrew_McCallum) ▾  
Andrew McCallum is a professor and researcher in the computer science department at University of Massachusetts Amherst. His primary specialties are in ...

[Andrew Mccallum - United Kingdom profiles | LinkedIn](#)  
[uk.linkedin.com/pub/dir/Andrew/Mccallum](http://uk.linkedin.com/pub/dir/Andrew/Mccallum) ▾  
View the profiles of professionals on LinkedIn named Andrew Mccallum located in the United Kingdom. There are 25 professionals named Andrew Mccallum in ...

**Andrew McCallum**  
Software Developer

Andrew McCallum is a professor and researcher in the computer science department at University of Massachusetts Amherst. [Wikipedia](#)



**Education:** Dartmouth College, University of Rochester  
**Awards:** Best 10-year Paper Award of the ICML

**People also search for**



[Tom M. Mitchell](#) [Lee Giles](#) [David M. Blei](#) [Michael Collins](#) [Robert Schapire](#)

[Feedback/More info](#)

# Machine Reading

**bing** Andrew McCallum | 

382,000 RESULTS Any time ▾

[Andrew McCallum Homepage - UMass CS | School of Computer ...](#)  
www.cs.umass.edu/~mccallum ▾  
Machine learning, text and information retrieval and extraction, reinforcement learning.

[Andrew McCallum - Wikipedia, the free encyclopedia](#)  
en.wikipedia.org/wiki/Andrew\_McCallum ▾  
Andrew McCallum is a professor and researcher in the computer science department at University of Massachusetts Amherst. His primary specialties are in machine learning, natural language processing, informa... + en.wikipedia.org

[Images of Andrew McCallum](#)  
bing.com/images



[Andrew McCallum - Twitter](#)  
www.twitter.com/andrew\_mccallum  
We would like to show you a description here but the site won't allow us.

[andrew mccallum profiles | LinkedIn](#)  
www.linkedin.com/pub/dir/andrew/mccallum ▾  
View the profiles of professionals named **andrew mccallum** on LinkedIn. There are 25 professionals named **andrew mccallum**, who use LinkedIn to exchange information ...

**Andrew McCallum**

Andrew McCallum is a professor and researcher in the computer science department at University of Massachusetts Amherst. His primary specialties are in machine learning, natural language processing, informa... + en.wikipedia.org

Education: Dartmouth College · University of Rochester

Awards: Best 10-year Paper Award of the ICML

---

People also search for

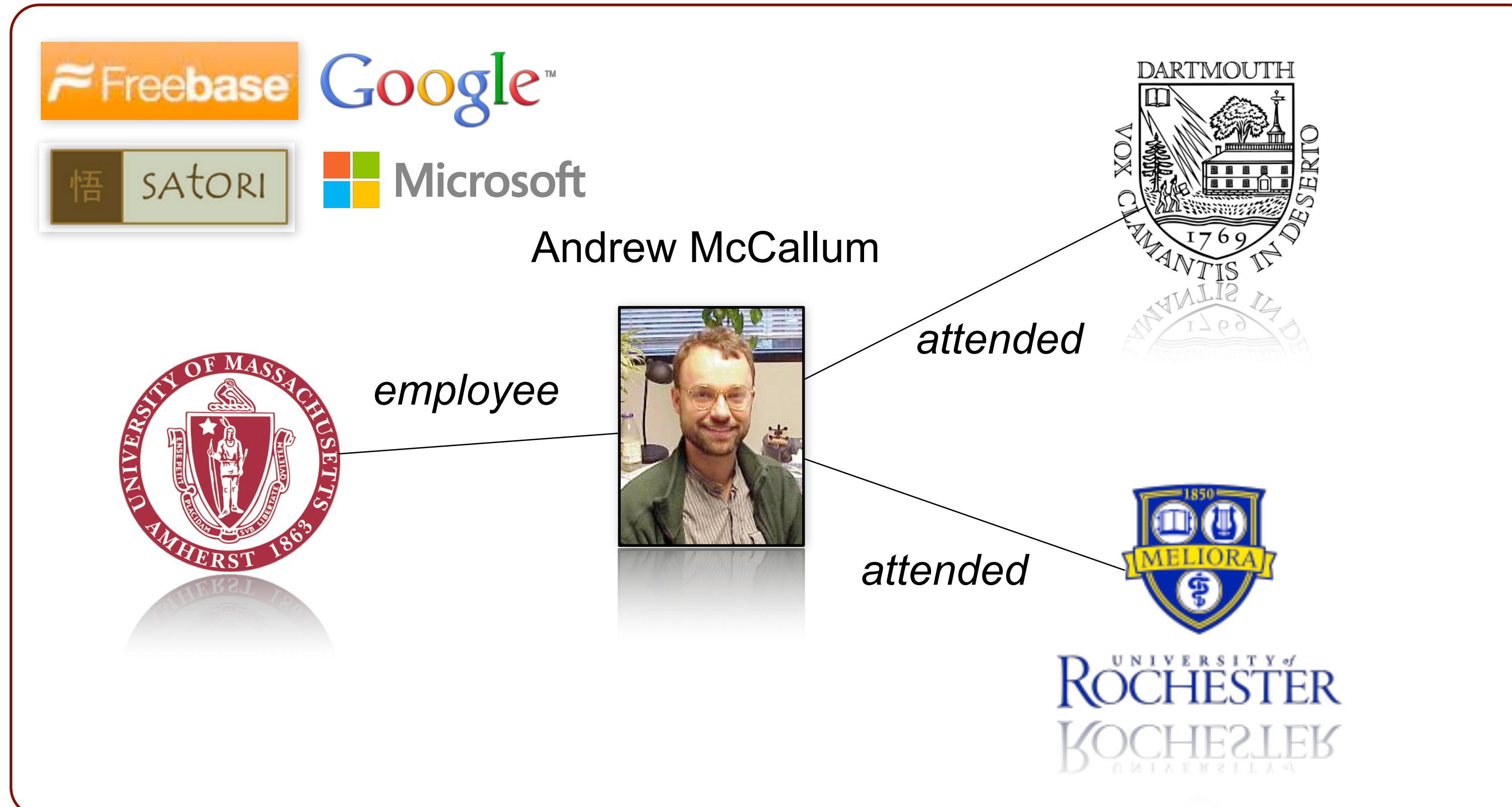


Tom M. Mitchell   Peter Norvig   Scott Fahlman

---

[Report a problem](#)

# Machine Reading



# Machine Reading



Professor  
**Andrew McCallum**  
Computer Science Department  
University of Massachusetts Amherst  
mccallum@cs.umass.edu  
+1 413 545-1323 (vox)  
+1 413 545-1789 (fax)



**Current Bio-sketch**

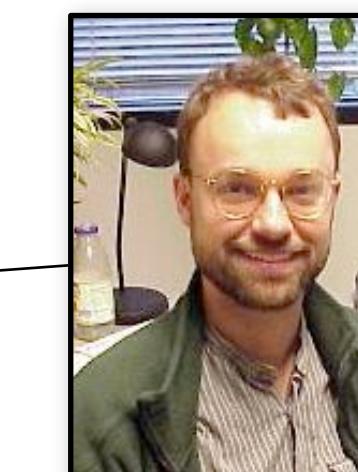
Andrew McCallum is a Professor and Director of the Information Extraction and Synthesis Laboratory in the School of Computer Science at University of Massachusetts Amherst. He has published over 250 papers in many areas of AI, including natural language processing, machine learning, data mining and reinforcement learning, and his work has received over 30,000 citations.

He obtained his PhD from University of Rochester in 1995 with Dana Ballard and a postdoctoral fellowship from CMU with Tom Mitchell and Sebastian Thrun. In the early 2000's he was Vice President of Research and Development at WhizBang Labs, a 170-person start-up company that used machine learning for information extraction from the Web.

**Contact**  
**Bio Vita**  
**Publications**  
**Talks**  
**Projects**



employee



attended



UNIVERSITY of  
**ROCHESTER**  
ROCHESTER

# Usages of Structured Knowledge

- Search
- Question Answering
- Data Mining
- Intelligent Agents
- Visualization

# Tasks

- Named Entity Recognition
- Relation Extraction
- Coreference
- Entity Linking
- Semantic Parsing
- Reasoning

# Tasks considered here

- Named Entity Recognition
- Relation Extraction
- Coreference
- Entity Linking
- Semantic Parsing
- Reasoning

# Named Entity Recognition

Identify mentions of entities in text (will become nodes in graph)

[Org] US President Obama spoke to ...  
[Per]  
...journalists in the White House  
[Loc]

# Named Entity Recognition

Needed in all kinds of domains

[Phos]

Phosphorylation of TRAF2 inhibits...

[Prot]

[Reg]

... binding to the CD40 domain

[Prot]

CD40

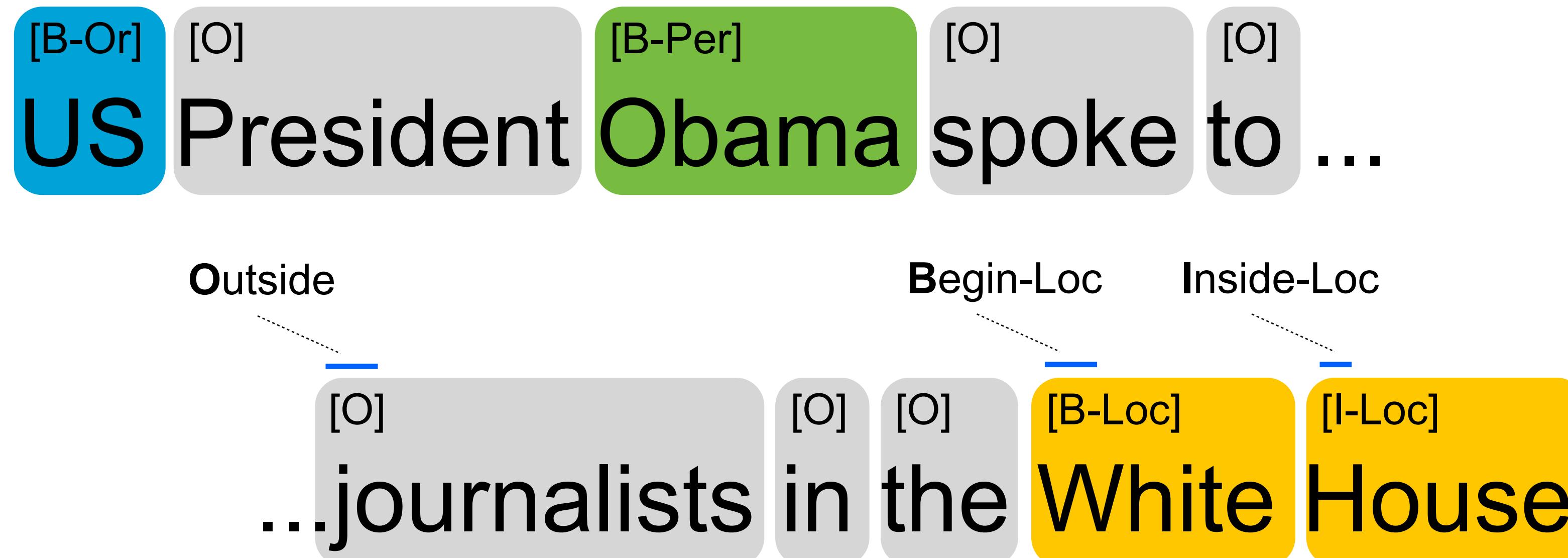
# Phrase Labeling

Label whole phrases with entity types (naively requires  $O(n^2)$  time)

[Org] US President Obama spoke to ...  
[Per]  
...journalists in the White House  
[Loc]

# Token Labeling (using BIO scheme)

Adapt label set and label on per-token basis (needs  $O(n)$ ) time



# Token Labeling

Useful for all kinds of things: Shallow Parsing (Chunking)

[B-NP] [I-NP] [I-NP]  
US President Obama spoke to ...

...journalists in the White House

# Token Labeling

Useful for all kinds of things: Part-of-Speech Tagging

[NNP] [NNP] [NNP] [VBD] [IN]  
US President Obama spoke to ...

...journalists in the White House

# Token Labeling

Useful for all kinds of things: Word Segmentation

[B] [B] [B] [I] [B] [B] [I]  
彼 は 音 楽 を 聞 く の が 大 好 き で す

# NER Evaluation

Precision:

$$\frac{\text{True positives (correctly predicted spans)}}{\text{All predicted spans}}$$

Recall:

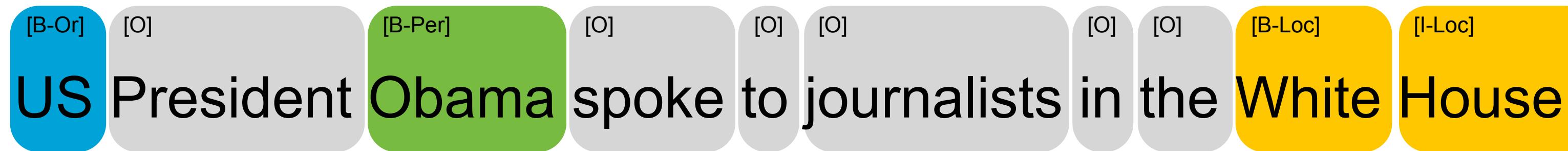
$$\frac{\text{True positives (correctly predicted spans)}}{\text{All true spans}}$$

F<sub>1</sub>:

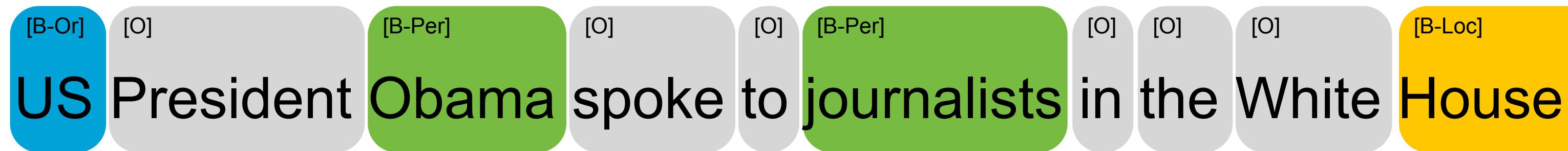
$$\frac{2 \cdot P \cdot R}{P + R} \quad \text{Harmonic Mean}$$

# NER Evaluation Example

TRUTH:



GUESS:

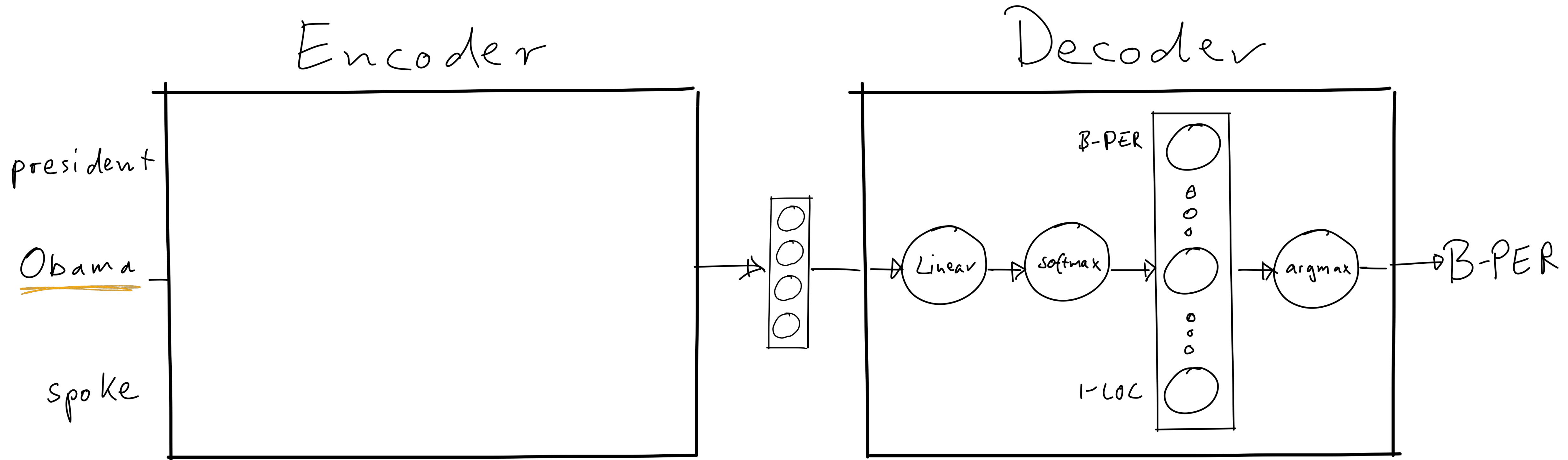


Precision:?

Recall:?

F1:?

# Local Model



# Feature Engineering

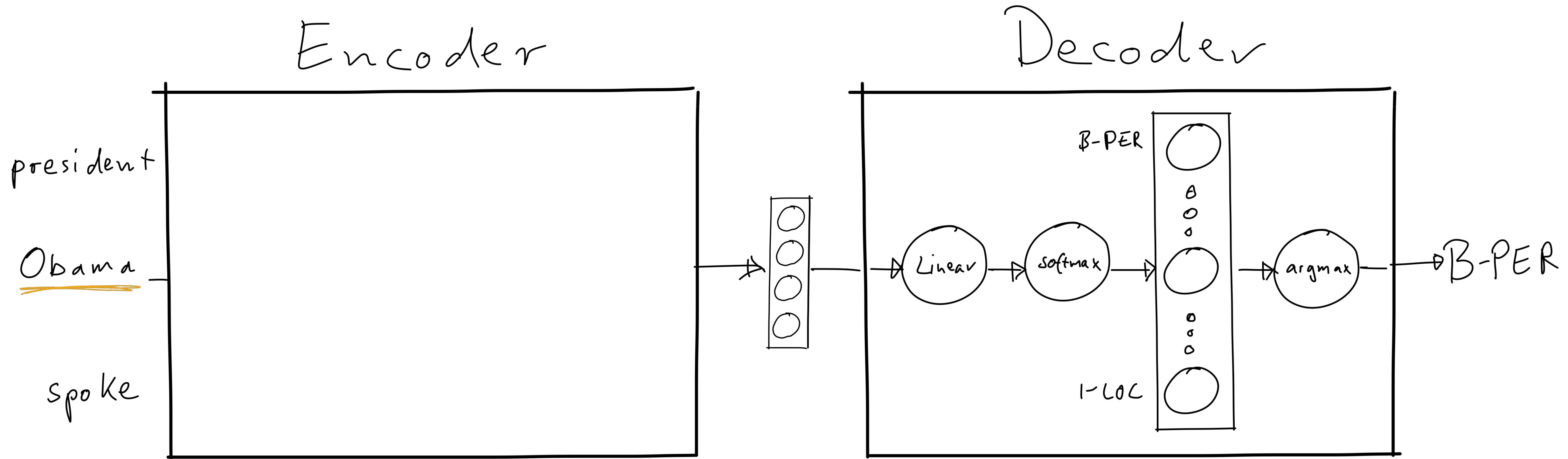
What could be good features (templates) for Named Entity Recognition?

[B-Or] [O] [B-Per] [O] [O]  
US President Obama spoke to ...

[O] [O] [O] [B-Loc] [I-Loc]  
...journalists in the White House

Bias Feature?

# Model 1



# Model 1

Only a *bias feature*

[O] [O] [O] [O] [O]  
US President Obama spoke to ...

[O] [O] [O] [O] [O]  
...journalists in the White House

(I am showing predictions on a development set after adding feature)

# Model 1

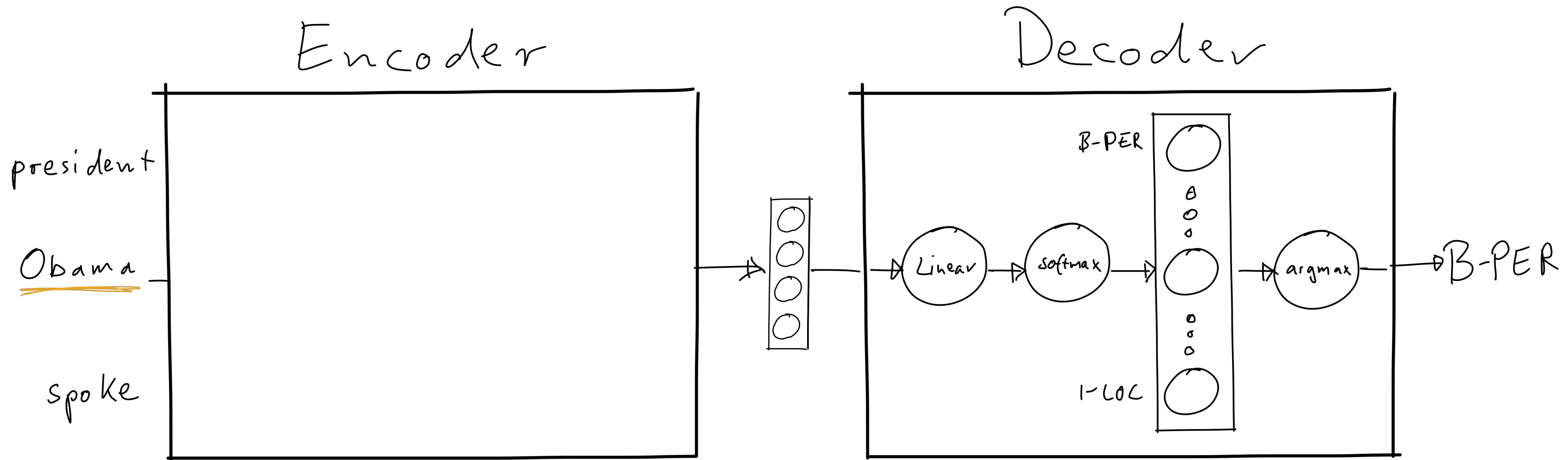
Only a bias

[O] [O]  
US President Obama spoke to ...

[O] [O] [O] [O]  
...journalists in the White House

Take into account words?

# Model 2



# Model 2

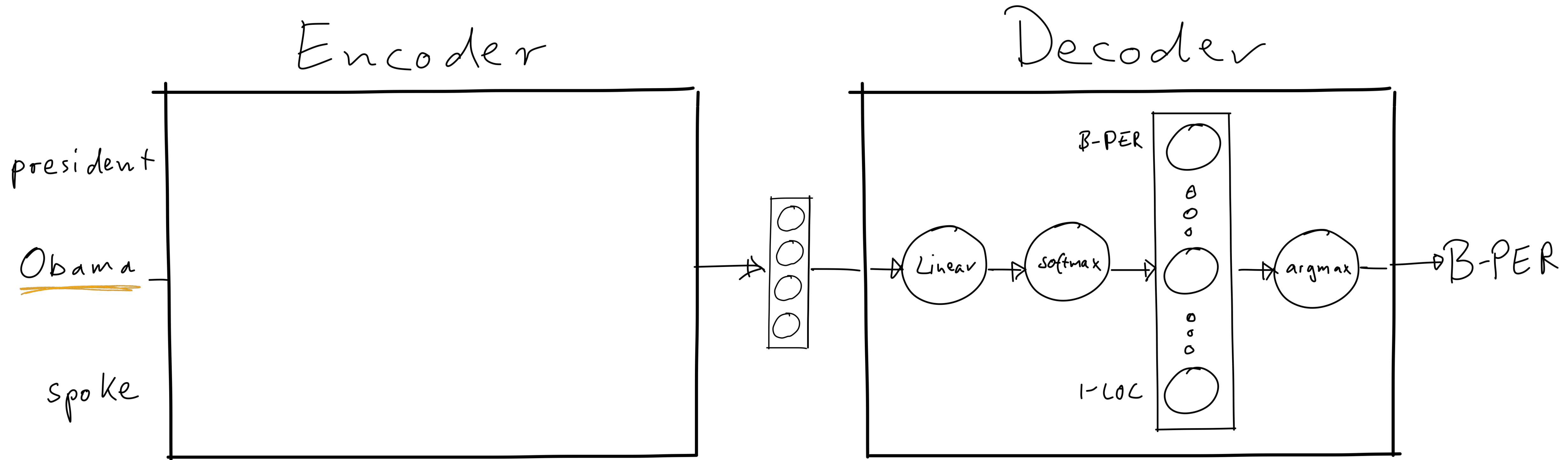
Add word identity

[B-Or] [O] [O] [O] [O]  
**US** President Obama spoke to ...

[O] [O] [O] [O] [O]  
...journalists in the White House

Take into account Capitalization?

# Model 3



# Model 3

With firstCapital

[B-Or] [O] [B-Loc] [O] [O]  
US President Obama spoke to ...

[O] [O] [O] [B-Loc] [B-Loc]  
...journalists in the White House

# Model 3

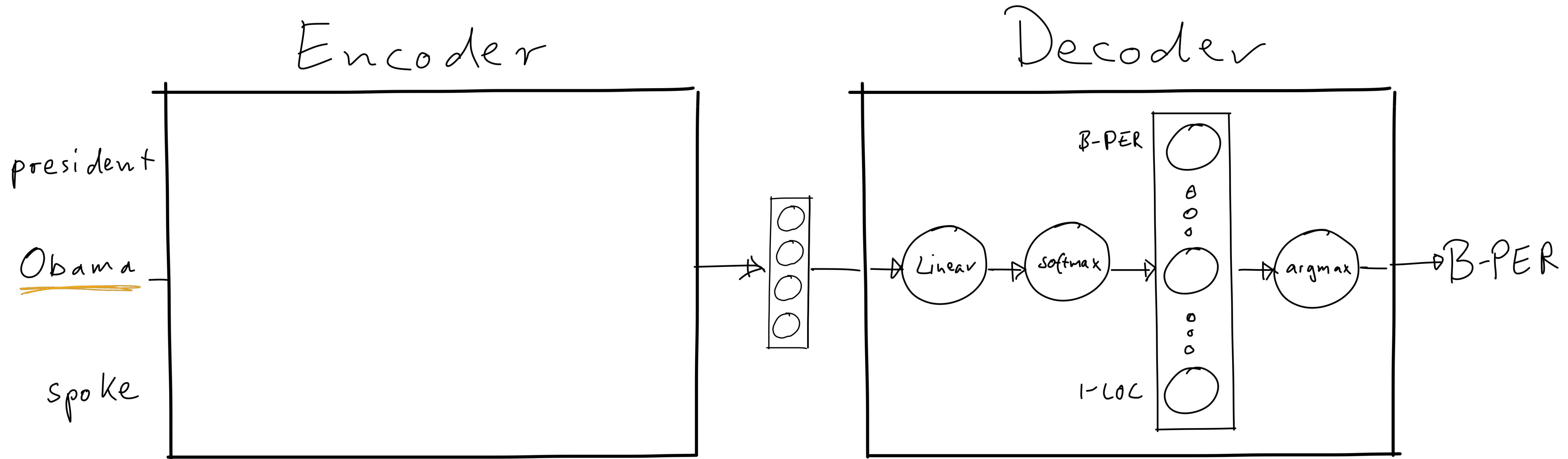
With firstCapital

[B-Or] [O] [B-Loc] [O] [O]  
US President Obama spoke to ...

[O] [O] [O] [B-Loc] [B-Loc]  
...journalists in the White House

Use a person last name dictionary feature

# Model 4



# Model 4

Use a person last name dictionary feature

[B-Or] [O] [B-Per] [O] [O]  
US President Obama spoke to ...

[O] [O] [O] [B-Loc] [B-Loc]  
...journalists in the White House

# Model 4

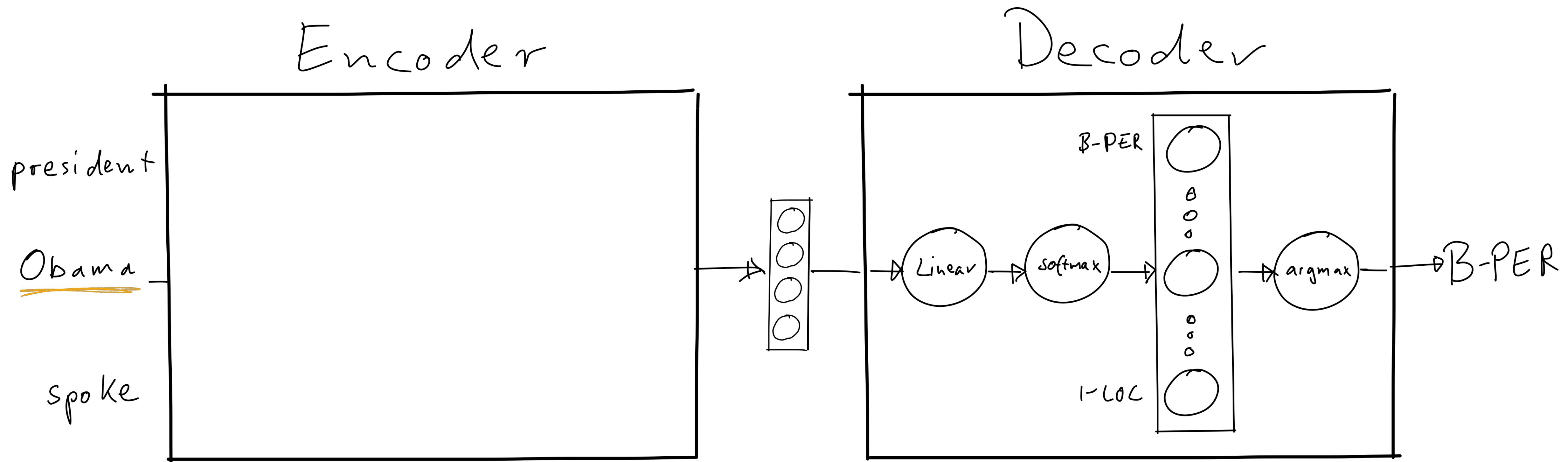
Use a person last name dictionary feature

[B-Or] [O] [B-Per] [O] [O]  
US President Obama spoke to ...

[O] [O] [O] [B-Loc] [B-Loc]  
...journalists in the White House

After capitalized words “I-” is more likely than “B-”

# Model 5



# Model 5

Previous word capitalized

[B-Or] [O] [I-Per] [O] [O]  
**US** President Obama spoke to ...

[O] [O] [O] [B-Loc] [I-Loc]  
...journalists in the White House

# Model 5

Previous word capitalized

[B-Or] [O] [I-Per]  
US President Obama spoke to ...

[O] [O] [O] [B-Loc] [I-Loc]  
...journalists in the White House

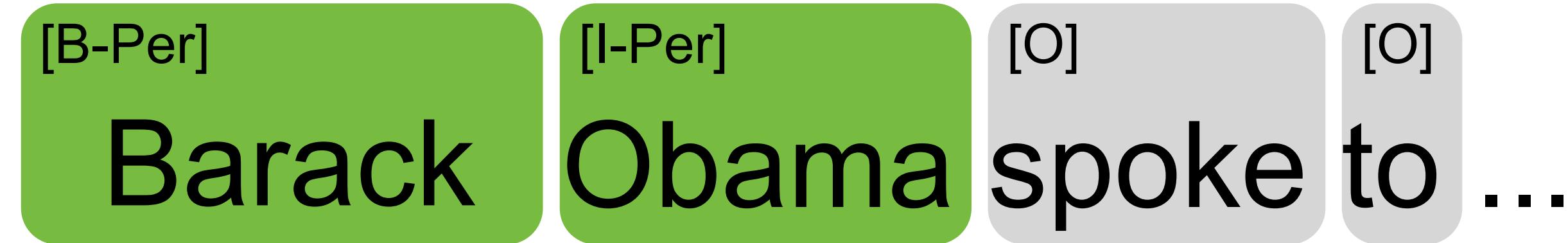
# Context Dependencies

Previous word affects the current label

[O] [I-Per] [O] [O]  
President Obama spoke to ...

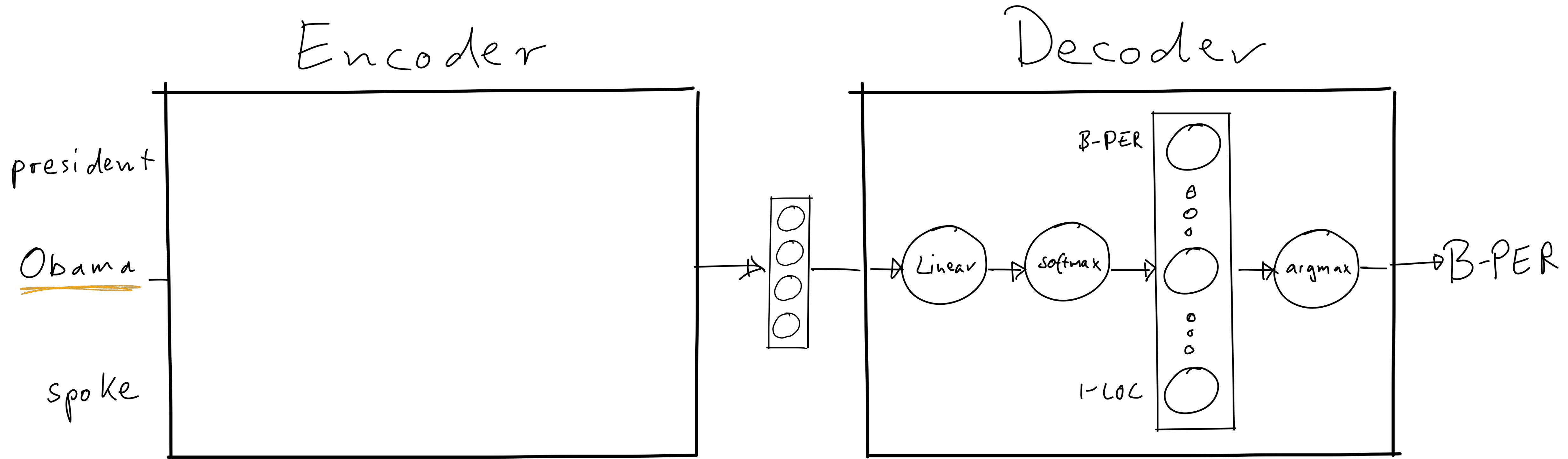
# Context Dependencies

Previous word affects the current label



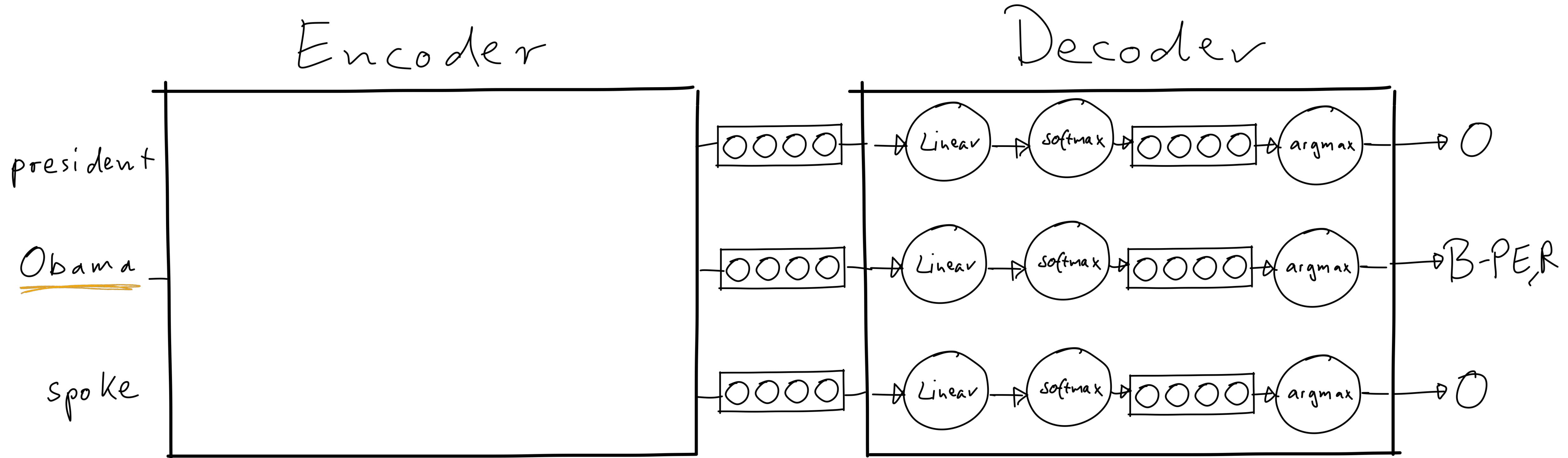
[B-Per]  
Barack [I-Per]  
Obama [O]  
[O] spoke to...

# Model 6



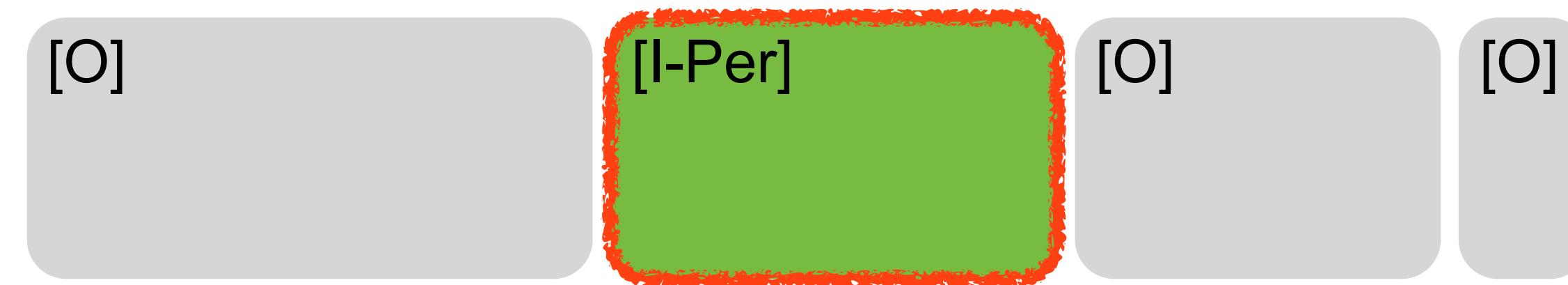
Complexity of RNN encoder  
per token?

# Reusing Computation



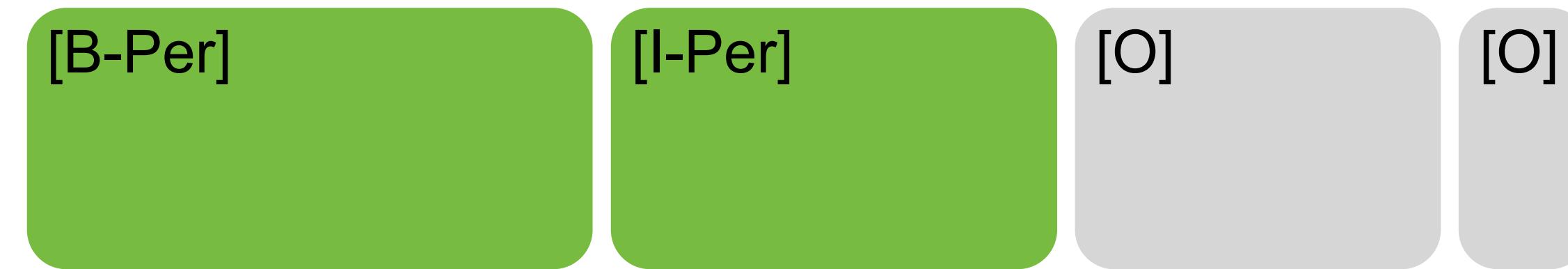
# Label Dependencies

Previous **label** affects the current label



# Label Dependencies

Previous **label** affects the current label



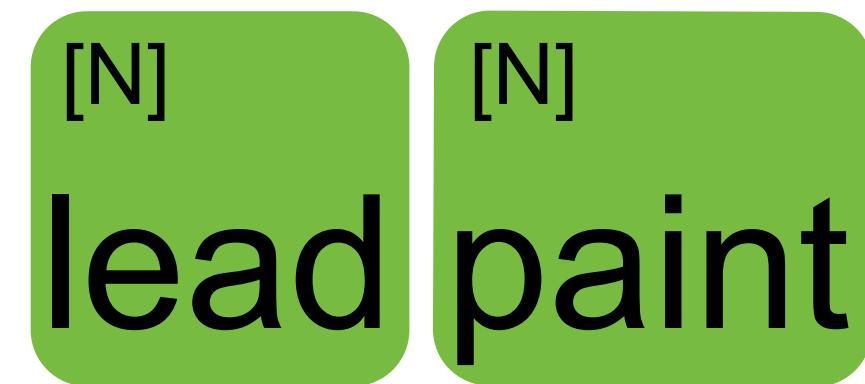
# Label Dependencies: Part-of-Speech

Previous label affects the current label



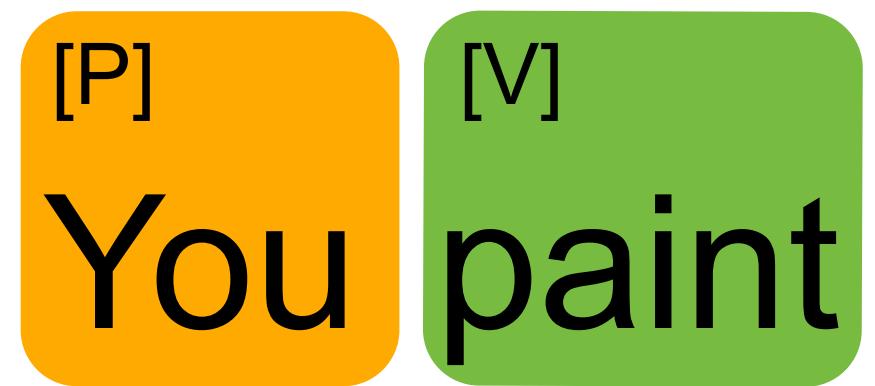
# Label Dependencies: Part-of-Speech

Previous label affects the current label

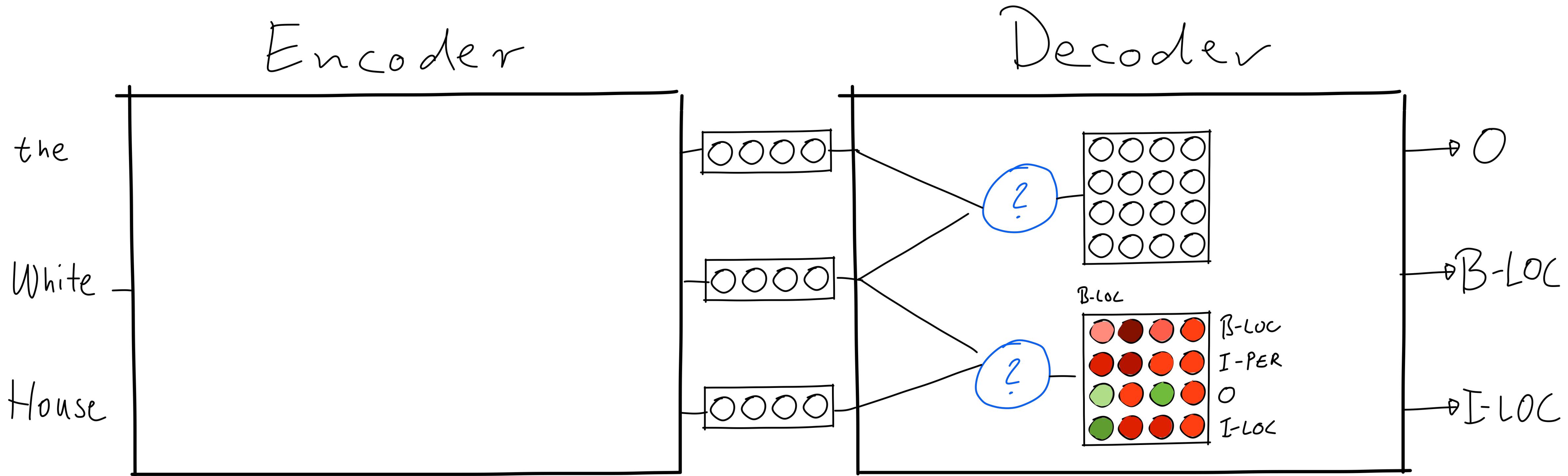


# Label Dependencies: Part-of-Speech

Previous label affects the current label



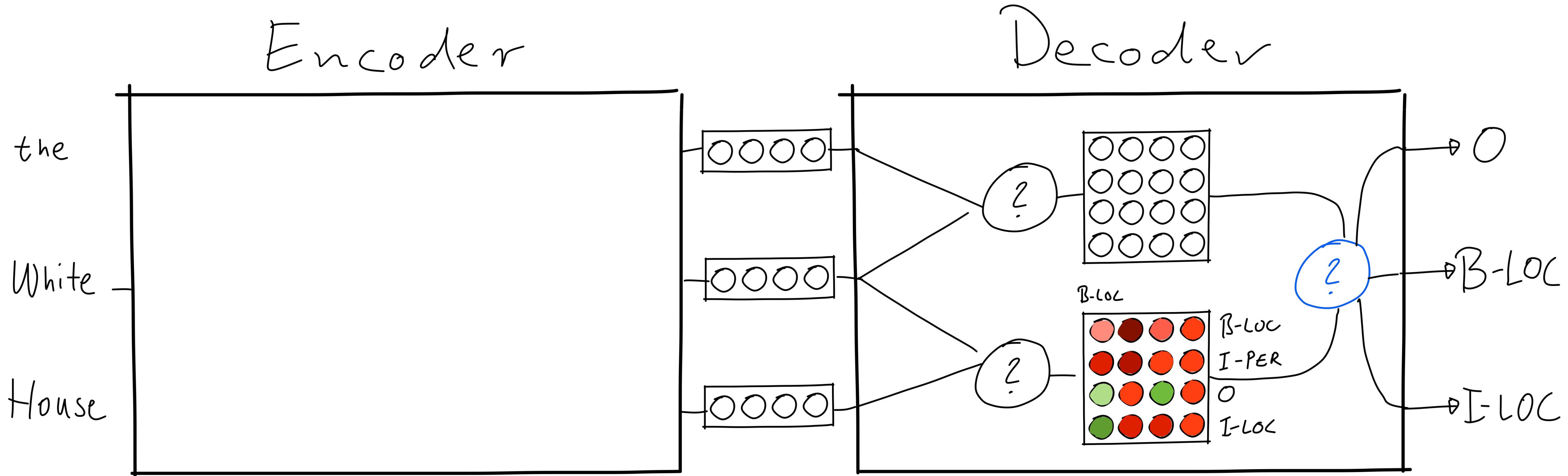
# Modelling Label Dependencies



$$s_i(y_i, y_{i+1}) = s_i(y_i) + s_{i+1}(y_{i+1}) + \theta_{y_i, y_{i+1}}$$

Local Score
Transition Matrix

# Modelling Label Dependencies



global score:  $s(y_1, \dots, y_n) = \sum_{i=1}^{n-1} s_i(y_i, y_{i+1})$

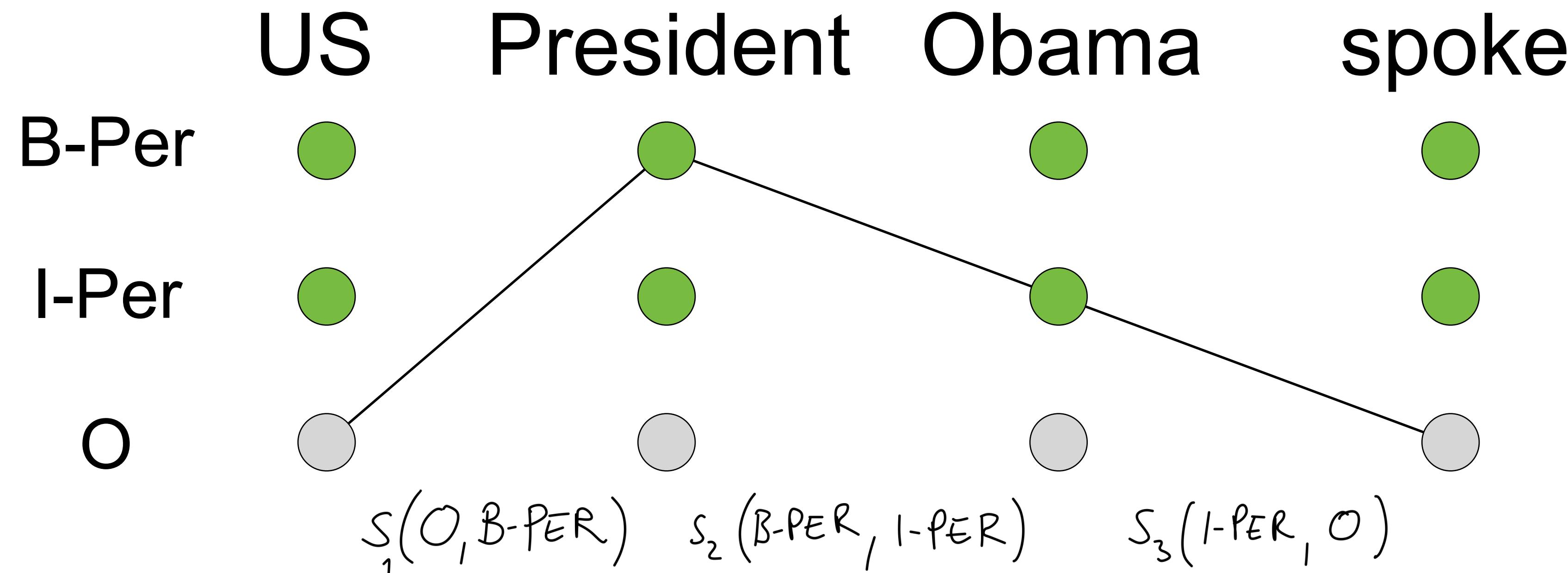
best label sequence:  $\hat{y} = \operatorname{argmax}_{y_1, \dots, y_n} s(y_1, \dots, y_n)$

Exhaustive

Search?

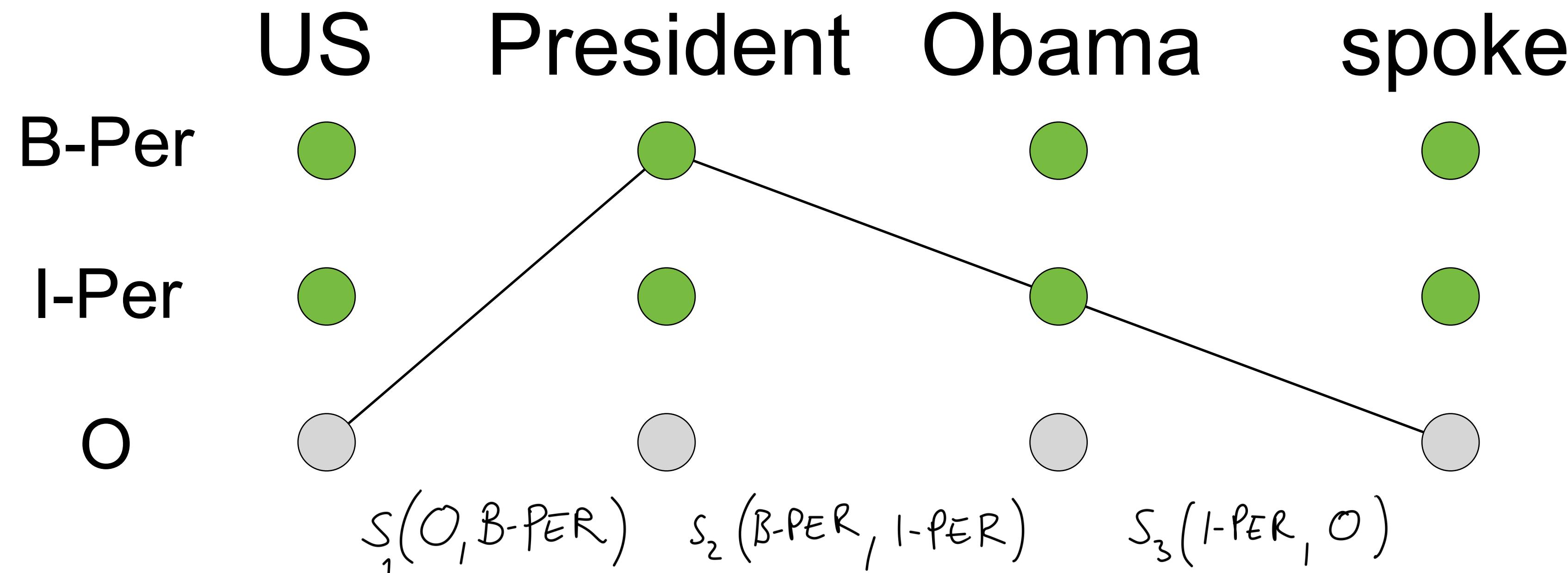
# Viterbi Decoding: Lattice

Representing a label sequence as a path, find highest scoring path



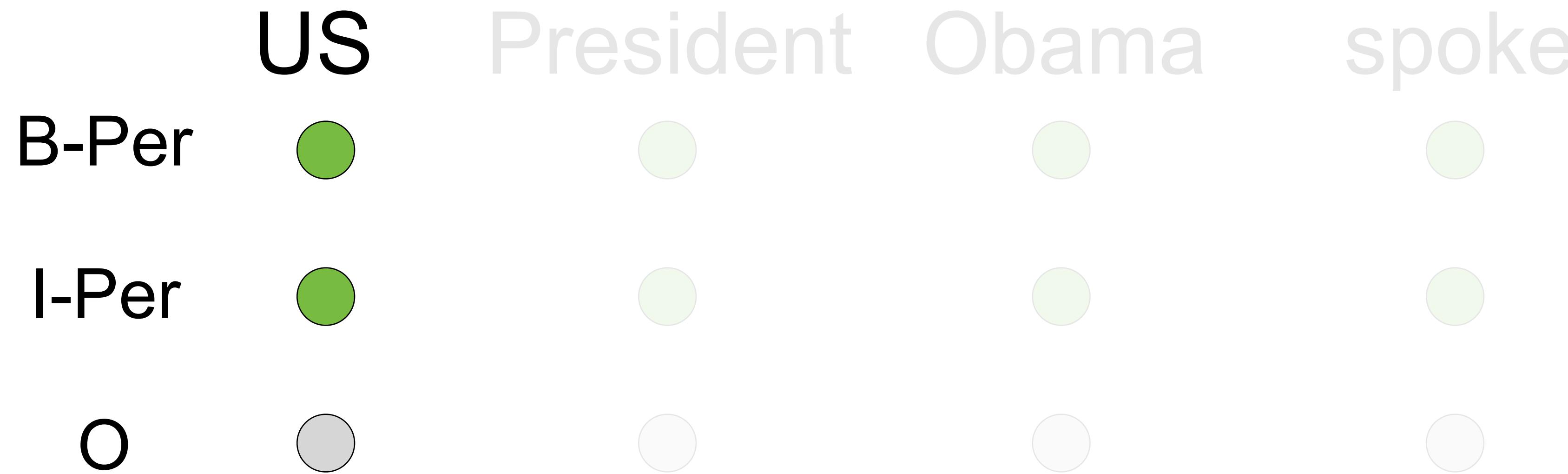
# Viterbi Decoding: Lattice

Representing a label sequence as a path, find highest scoring **path of length 4**



# Viterbi Decoding

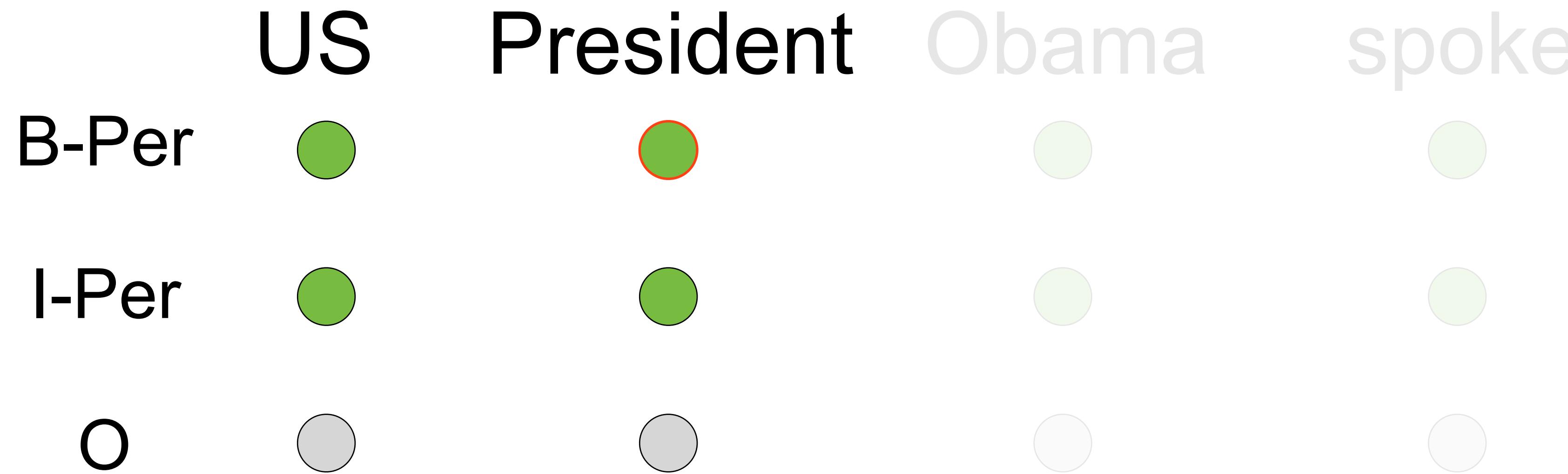
For each label, find best path of length 1 ending at this label, and its score



"Message": score of best sequence  
of length 1 ending in  $y_1$ .  $\alpha_1(y_1) = 0$

# Viterbi Decoding

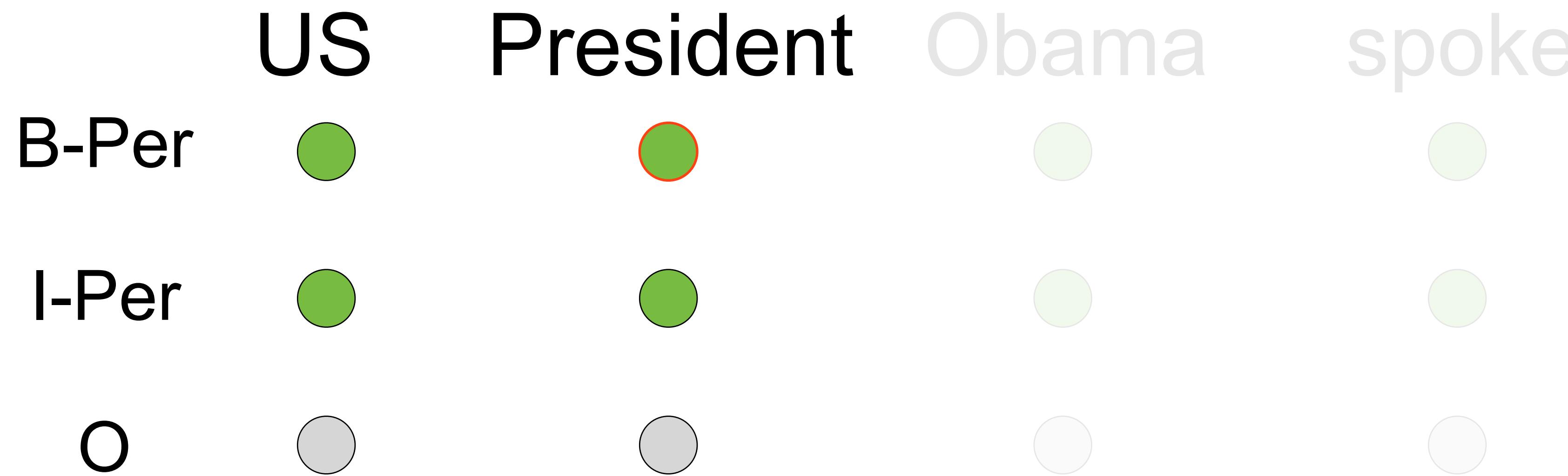
For each label, find best path of length 2 ending at this label, and its score



$$\alpha_2(y_2) = \max_{y_1} s_1(y_1, y_2) + \alpha_1(y_1)$$

# Viterbi Decoding

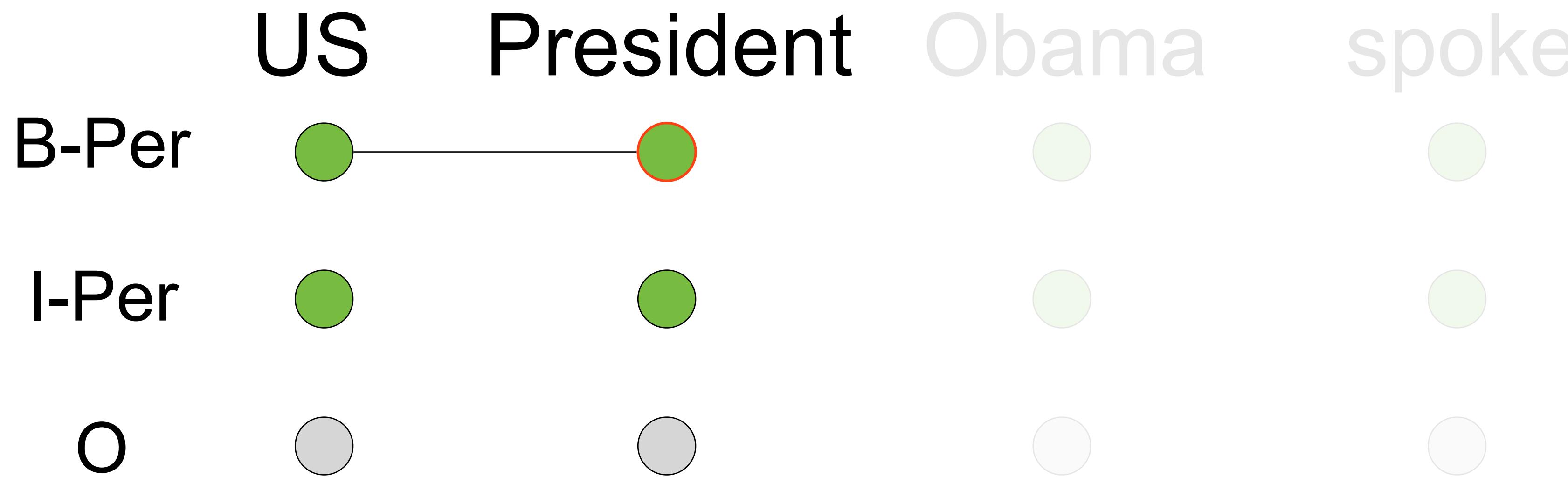
For each label, find best path of length 2 ending at this label, and its score



$$\alpha_2(y_2) = \max_{y_1} s_1(y_1, y_2) + \alpha_1(y_1)$$

# Viterbi Decoding

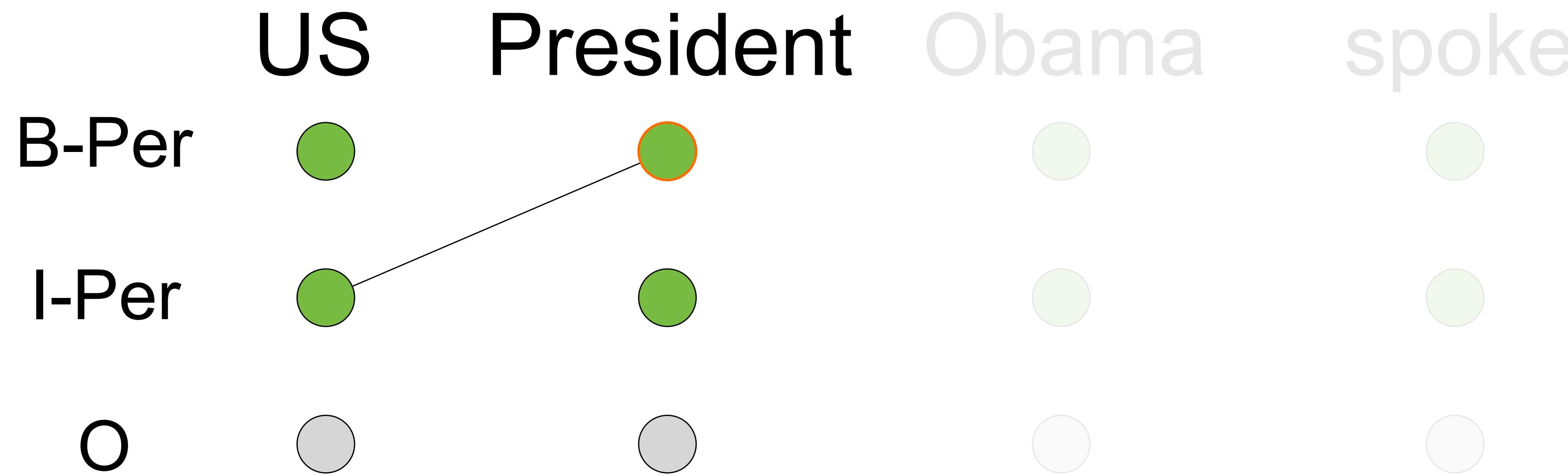
For each label, find best path of length 3 ending at this label, and its score



$$\alpha_2(y_2) = \max_{y_1} s_1(y_1, y_2) + \alpha_1(y_1)$$

# Viterbi Decoding

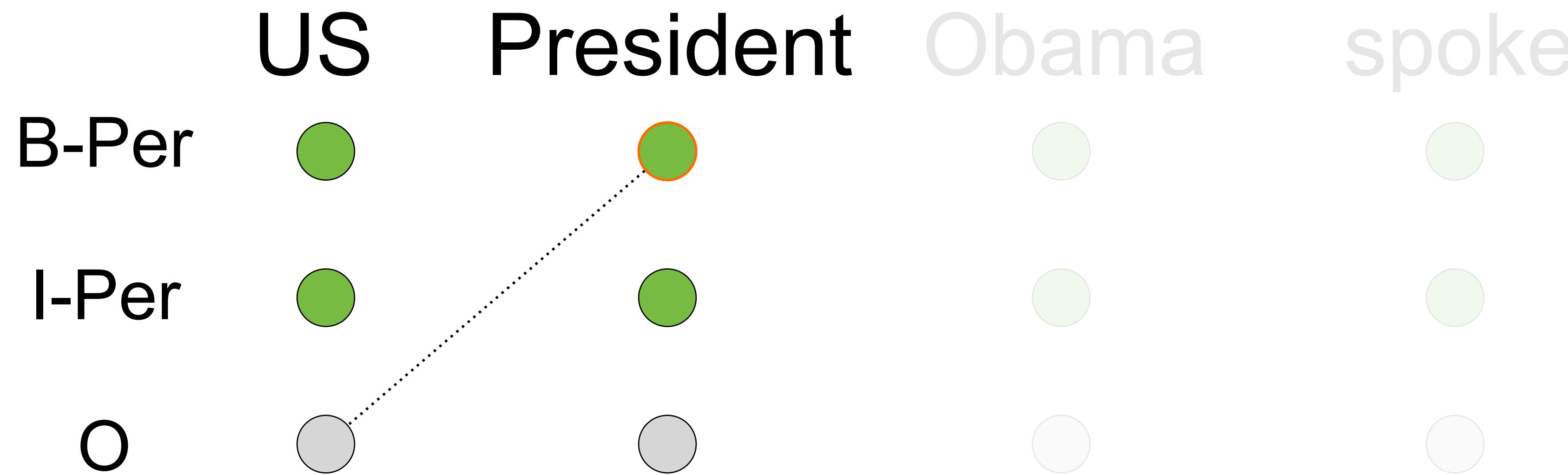
For each label, find best path of length 2 ending at this label, and its score



$$\alpha_2(y_2) = \max_{y_1} s_1(y_1, y_2) + \alpha_1(y_1)$$

# Viterbi Decoding

For each label, find best path of length 2 ending at this label, and its score

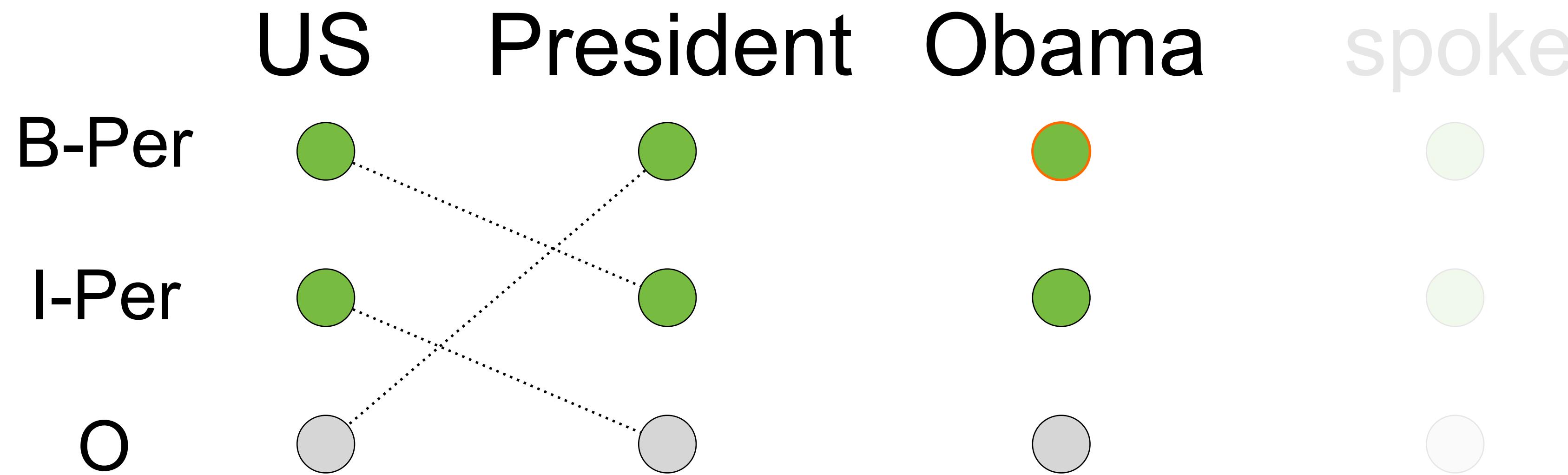


Remember  
winning  $y_1$

$$\alpha_2(y_2) = \max_{y_1} s_1(y_1, y_2) + \alpha_1(y_1)$$

# Viterbi Decoding

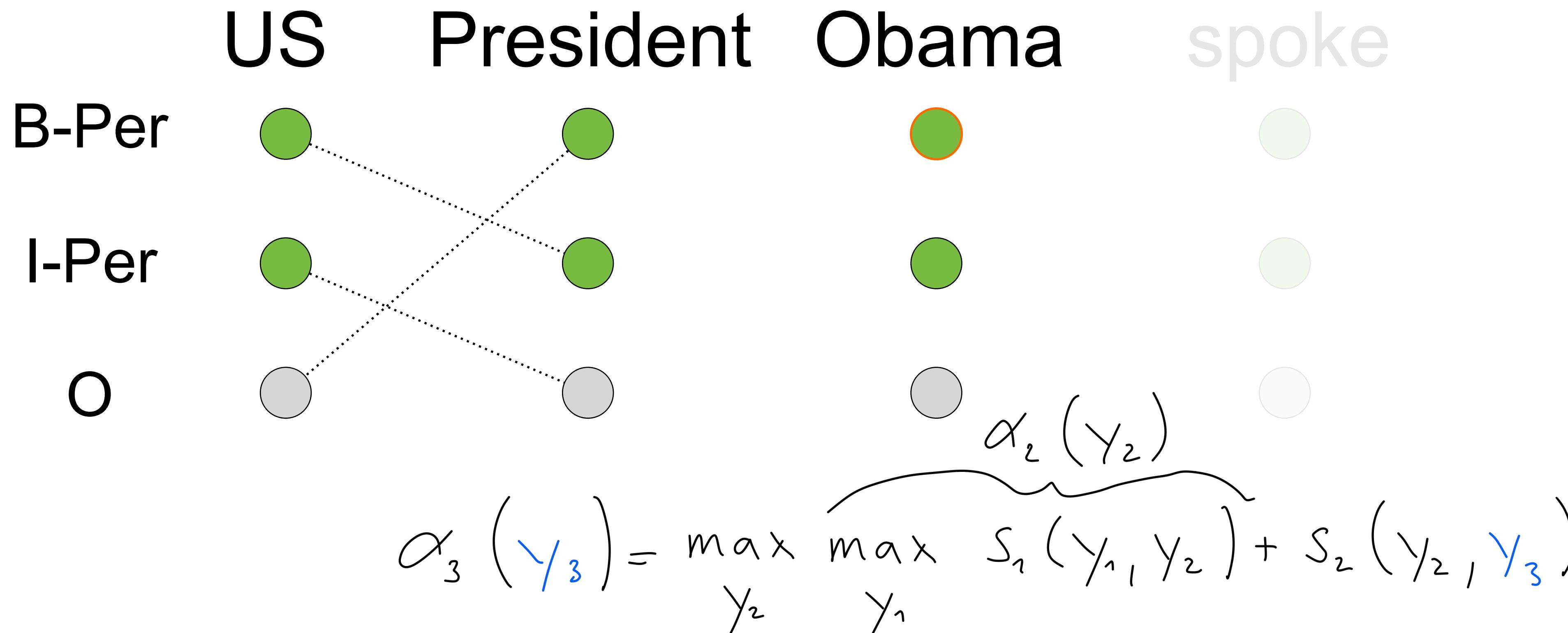
For each label, find best path of length 3 ending at this label, and its score



$$\alpha_3(y_3) = \max_{y_1, y_2} s_1(y_1, y_2) + s_2(y_2, y_3)$$

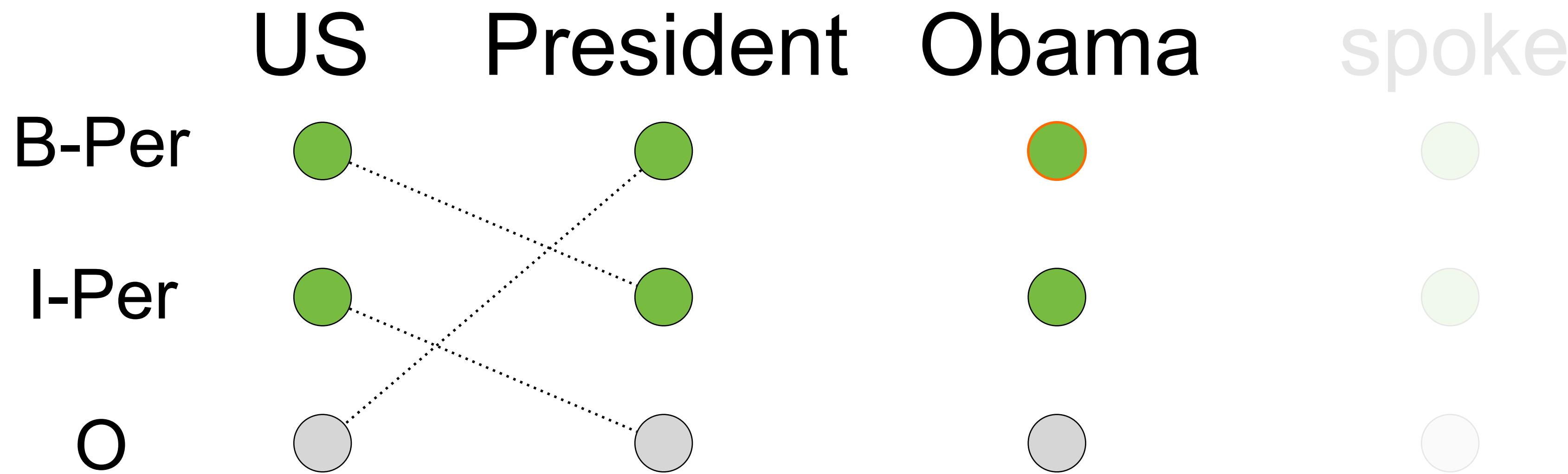
# Viterbi Decoding

For each label, find best path of length 3 ending at this label, and its score



# Viterbi Decoding

For each label, find best path of length 3 ending at this label, and its score



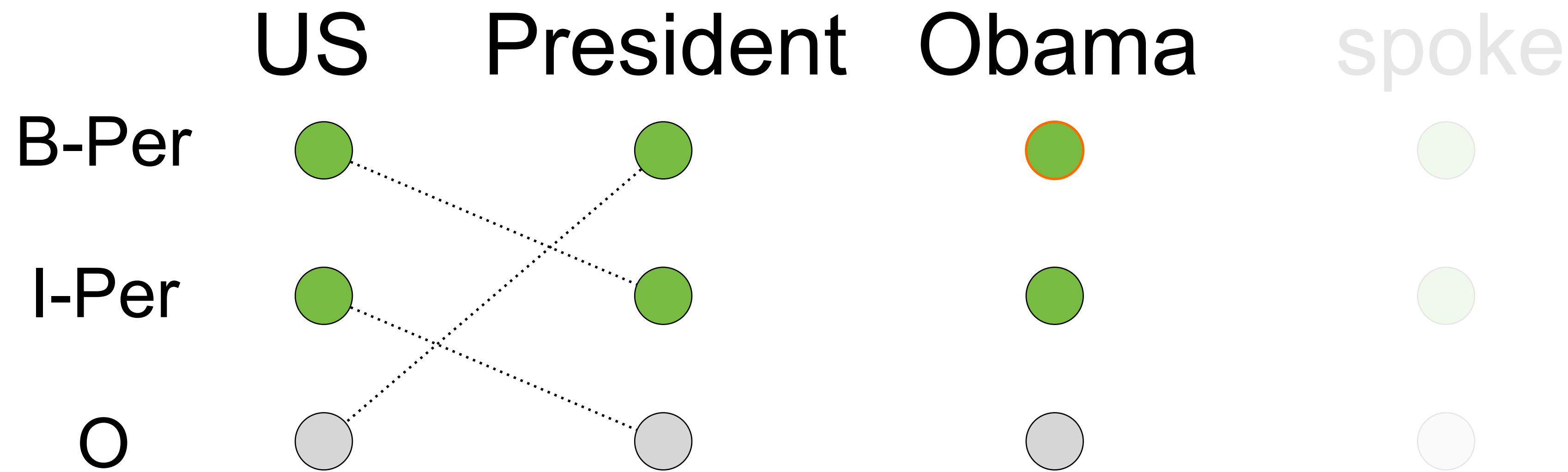
$$\alpha_3(y_3) = \max_{y_2} \quad$$

$$\alpha_2(y_2)$$

$$+ s_2(y_2, y_3)$$

# Viterbi Decoding

For each label, find best path of length 3 ending at this label, and its score



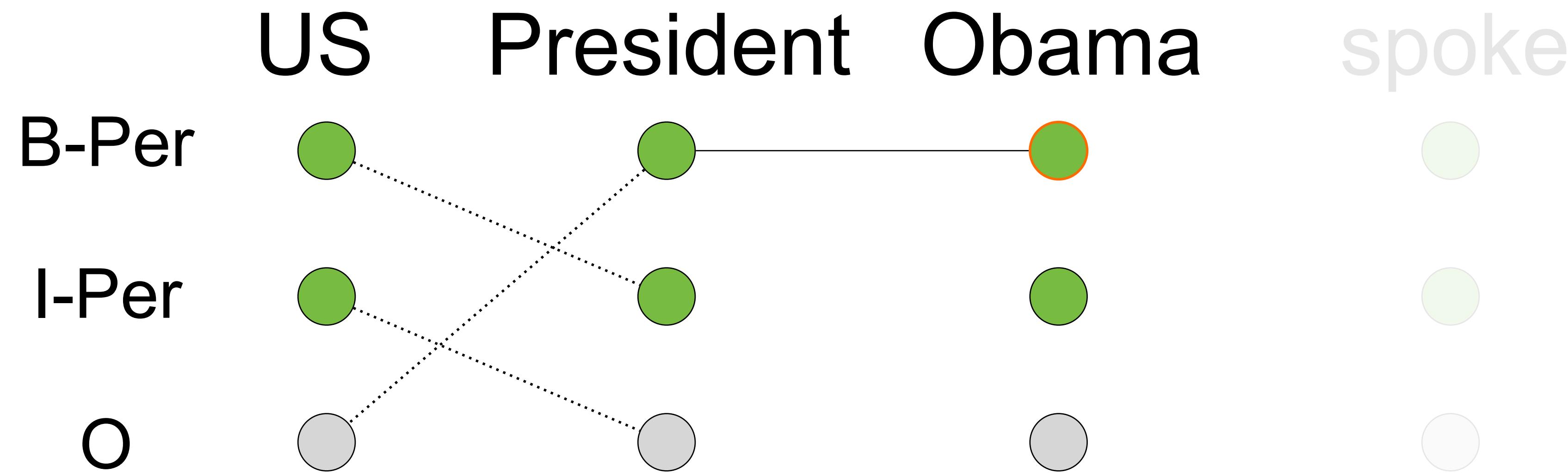
$$\alpha_3(y_3) = \max_{y_2} \quad$$

$$\alpha_2(y_2)$$

$$+ s_2(y_2, y_3)$$

# Viterbi Decoding

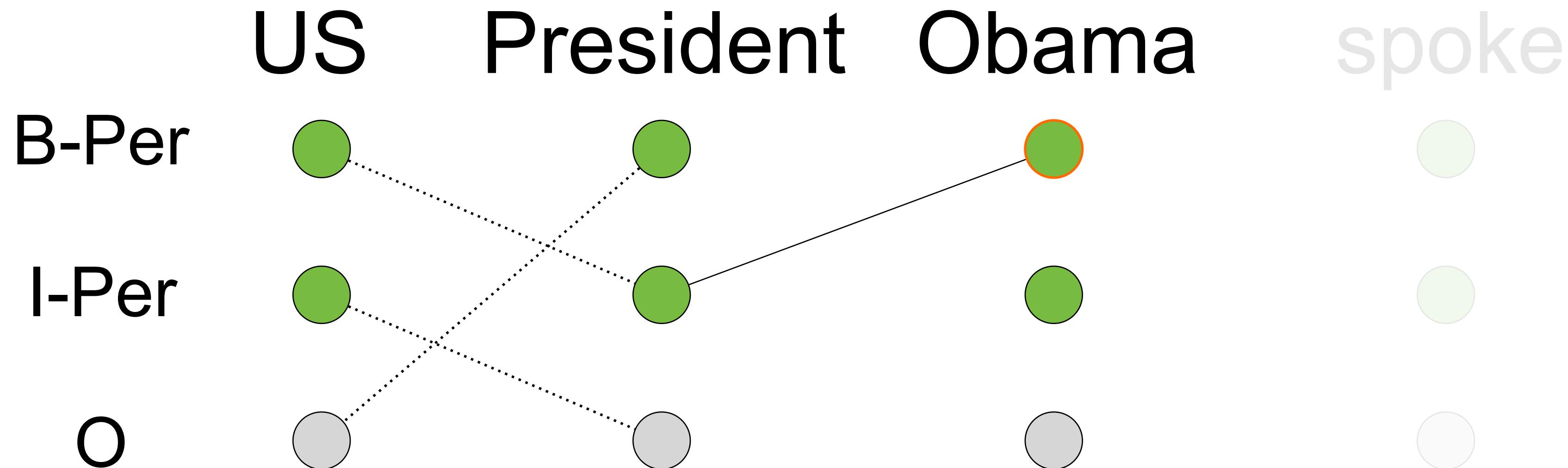
For each label, find best path of length 3 ending at this label, and its score



$$\alpha_3(y_3) = \max_{y_2} \alpha_2(y_2) + s_2(y_2, y_3)$$

# Viterbi Decoding

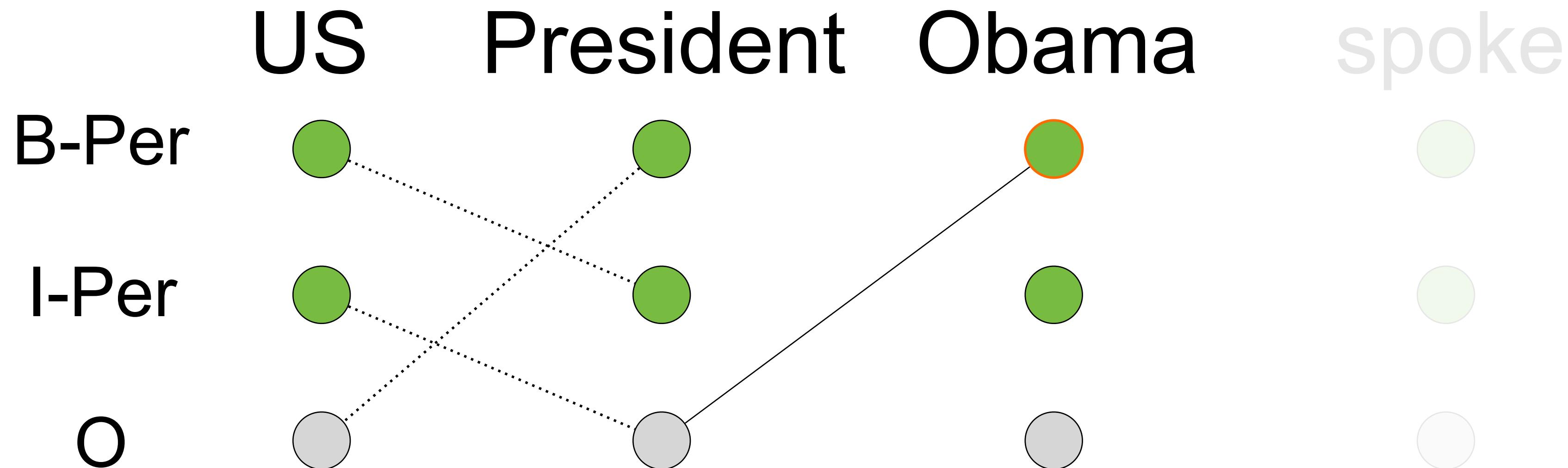
For each label, find best path of length 3 ending at this label, and its score



$$\alpha_3(y_3) = \max_{y_2} \alpha_2(y_2) + s_2(y_2, y_3)$$

# Viterbi Decoding

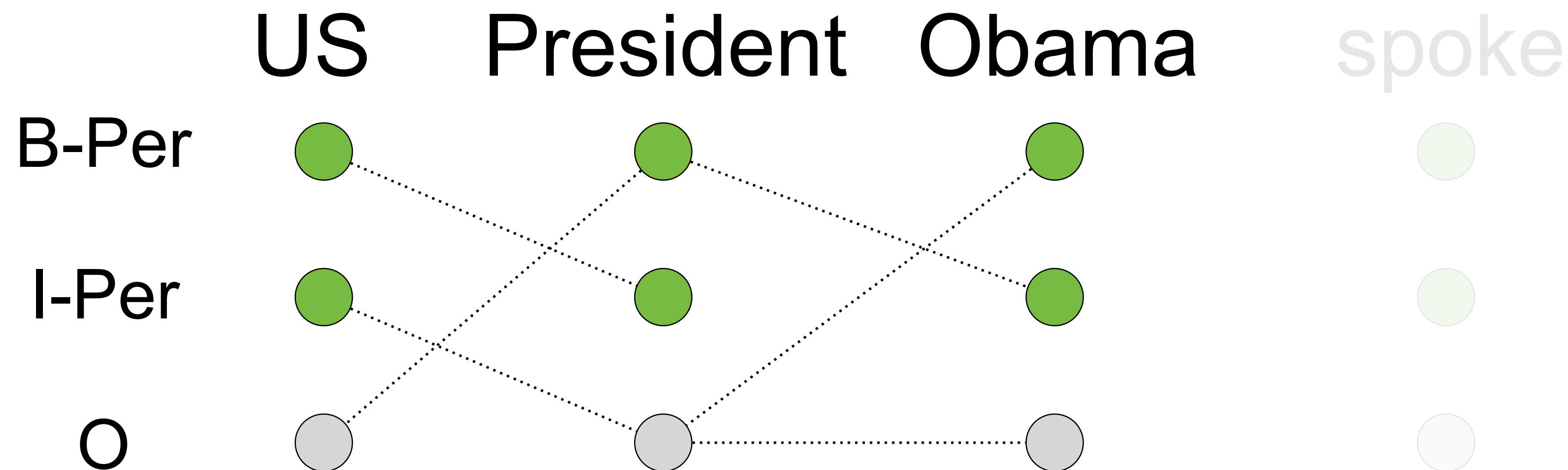
For each label, find best path of length 3 ending at this label, and its score



$$\alpha_3(y_3) = \max_{y_2} \alpha_2(y_2) + s_2(y_2, y_3)$$

# Viterbi Decoding

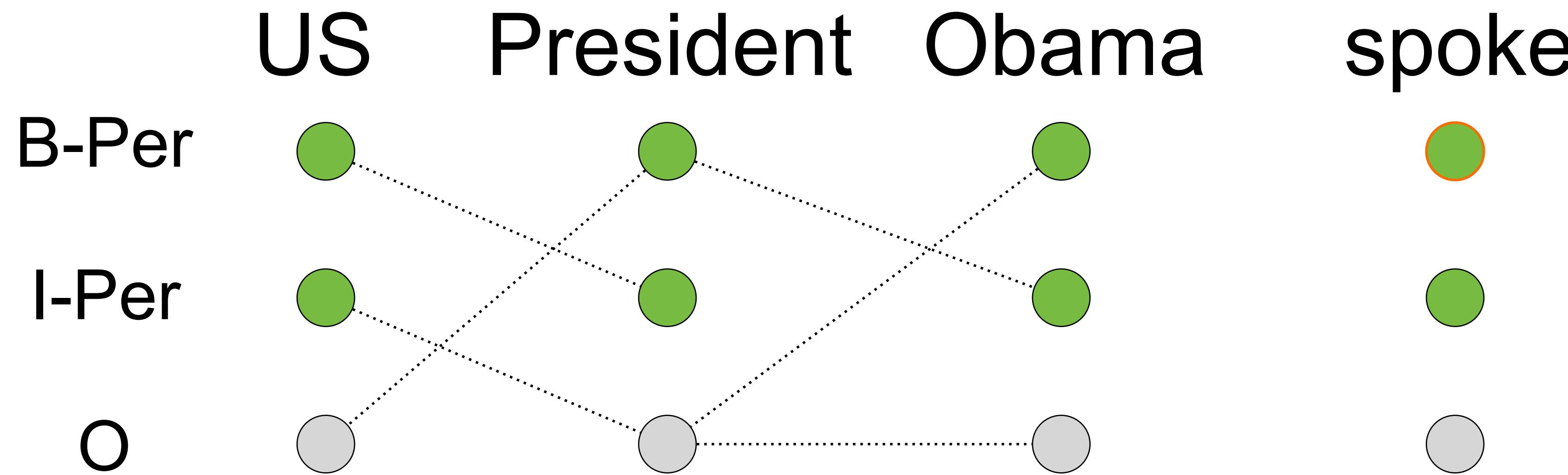
For each label, find best path of length 3 ending at this label, and its score



$$\alpha_3(y_3) = \max_{y_2} \alpha_2(y_2) + s_2(y_2, y_3)$$

# Viterbi Decoding

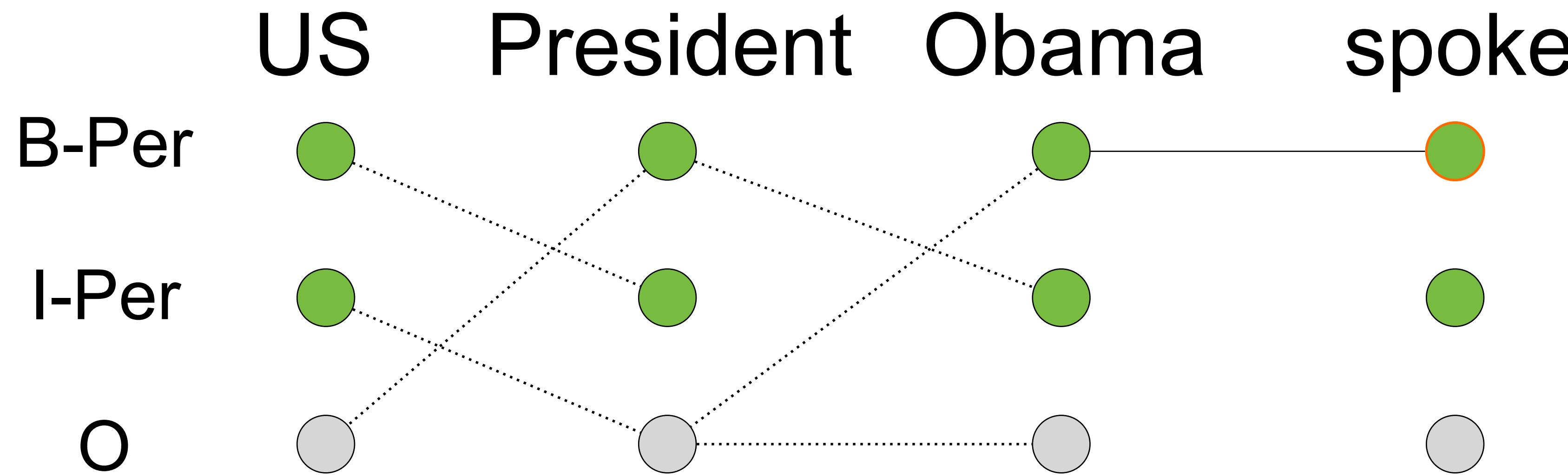
For each label, find best path of length 4 ending at this label, and its score



$$\alpha_4(y_4) = \max_{y_3} \alpha_3(y_3) + s_3(y_3, y_4)$$

# Viterbi Decoding

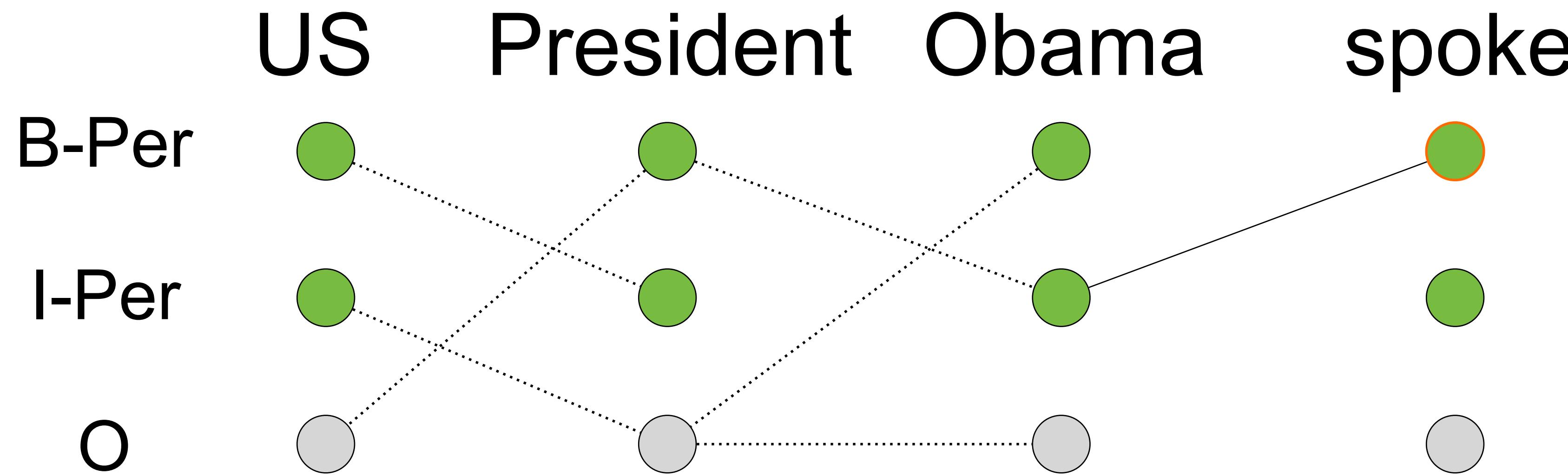
For each label, find best path of length 4 ending at this label, and its score



$$\alpha_4(y_4) = \max_{y_3} \alpha_3(y_3) + s_3(y_3, y_4)$$

# Viterbi Decoding

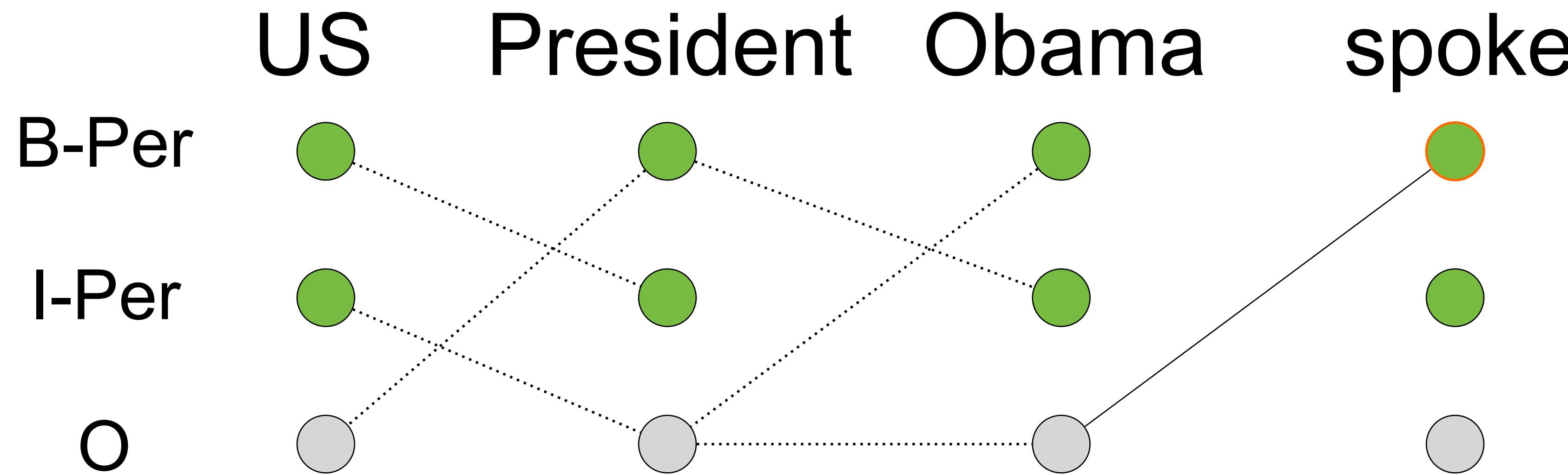
For each label, find best path of length 4 ending at this label, and its score



$$\alpha_4(y_4) = \max_{y_3} \alpha_3(y_3) + s_3(y_3, y_4)$$

# Viterbi Decoding

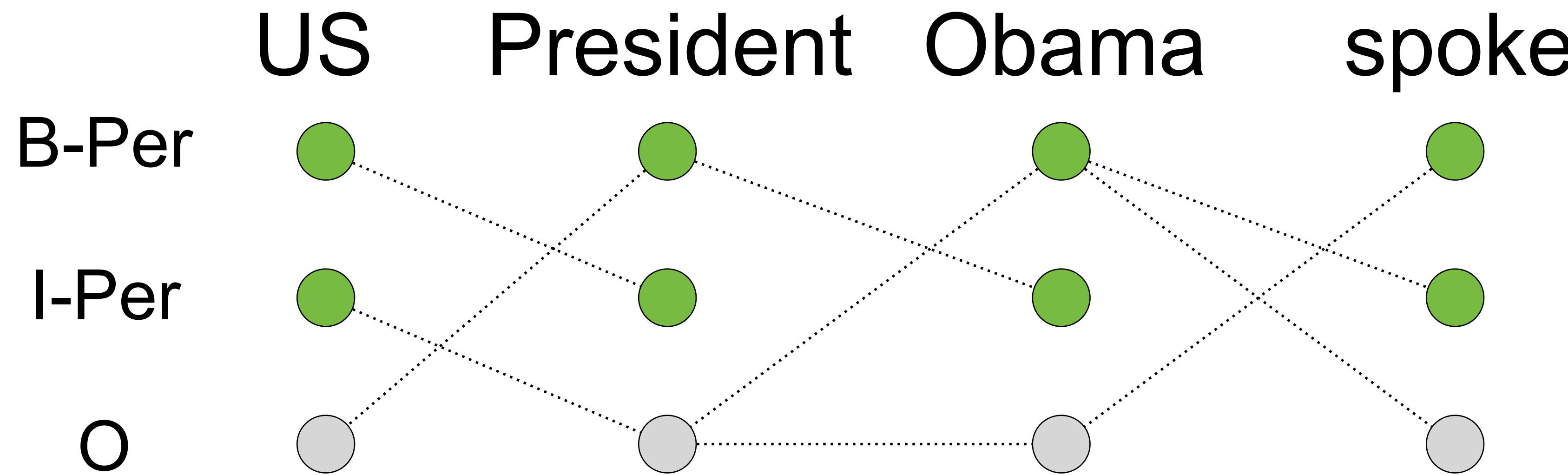
For each label, find best path of length 4 ending at this label, and its score



$$\alpha_4(y_4) = \max_{y_3} \alpha_3(y_3) + s_3(y_3, y_4)$$

# Viterbi Decoding

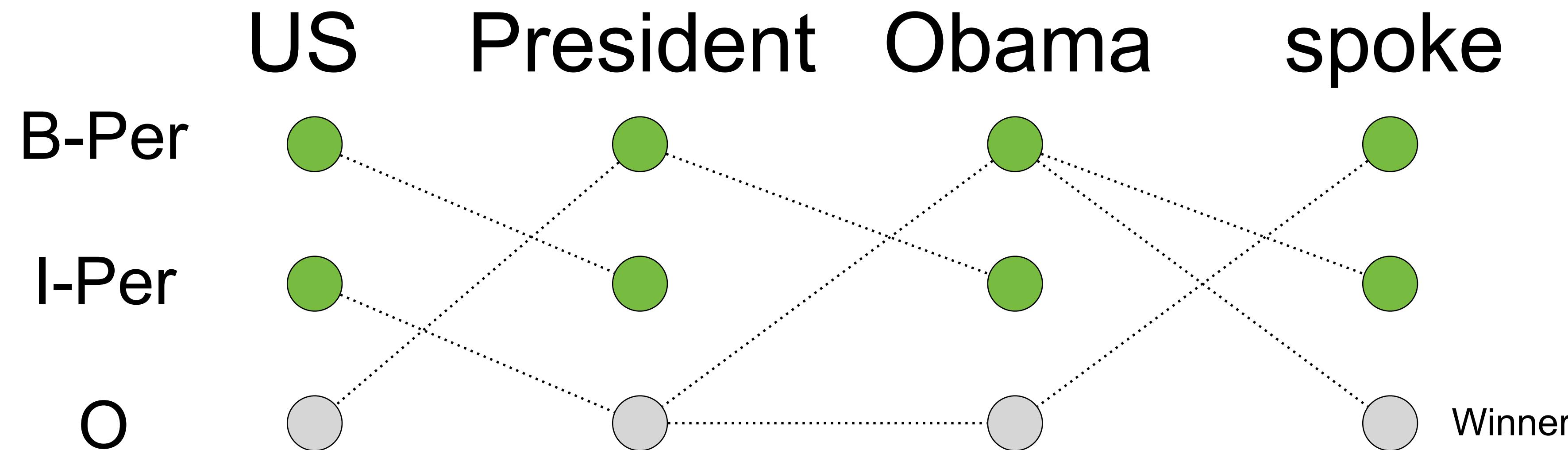
For each label, find best path of length 4 ending at this label, and its score



$$\alpha_4(y_4) = \max_{y_3} \alpha_3(y_3) + s_3(y_3, y_4)$$

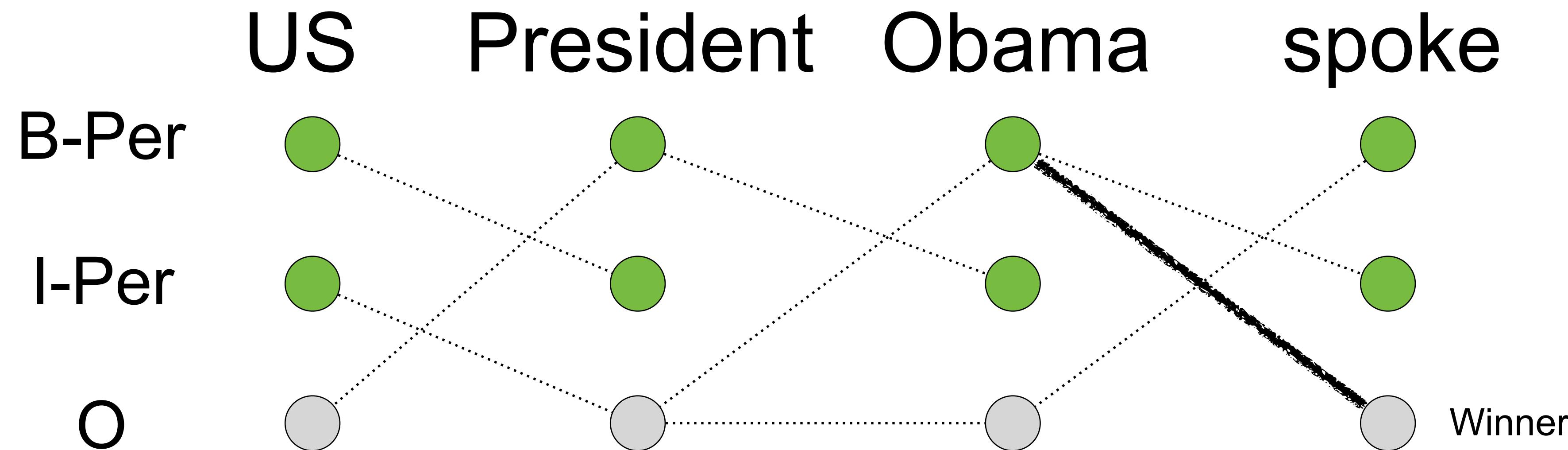
# Viterbi Decoding: Backtracking

You got a winner! Find the path that lead to it...



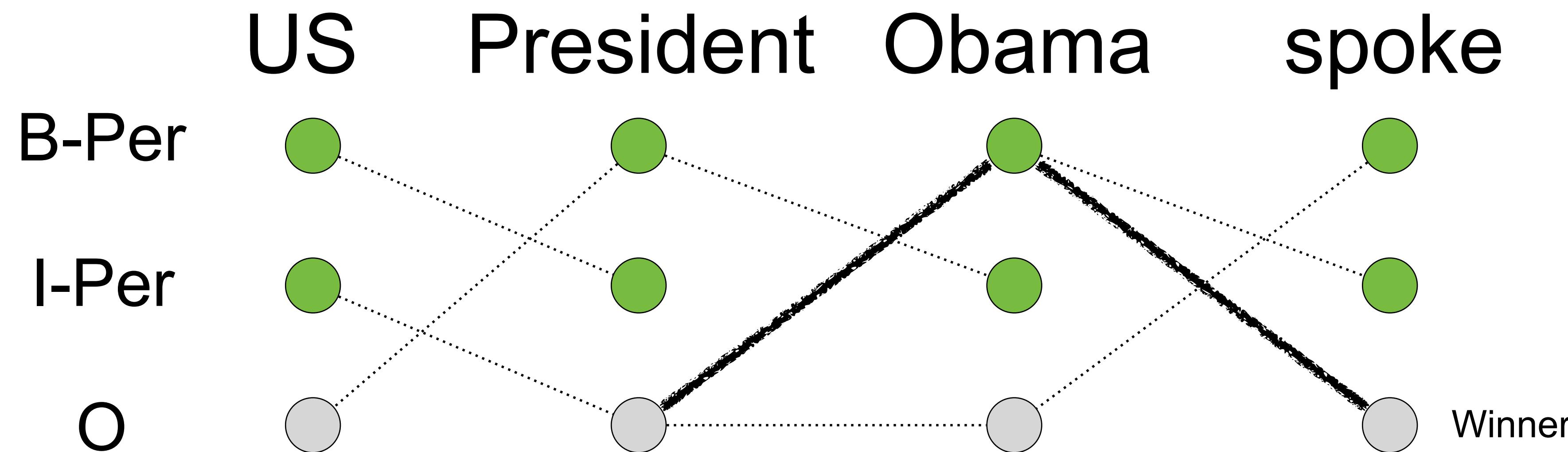
# Viterbi Decoding: Backtracking

Which was the maximizing node coming into label 4?



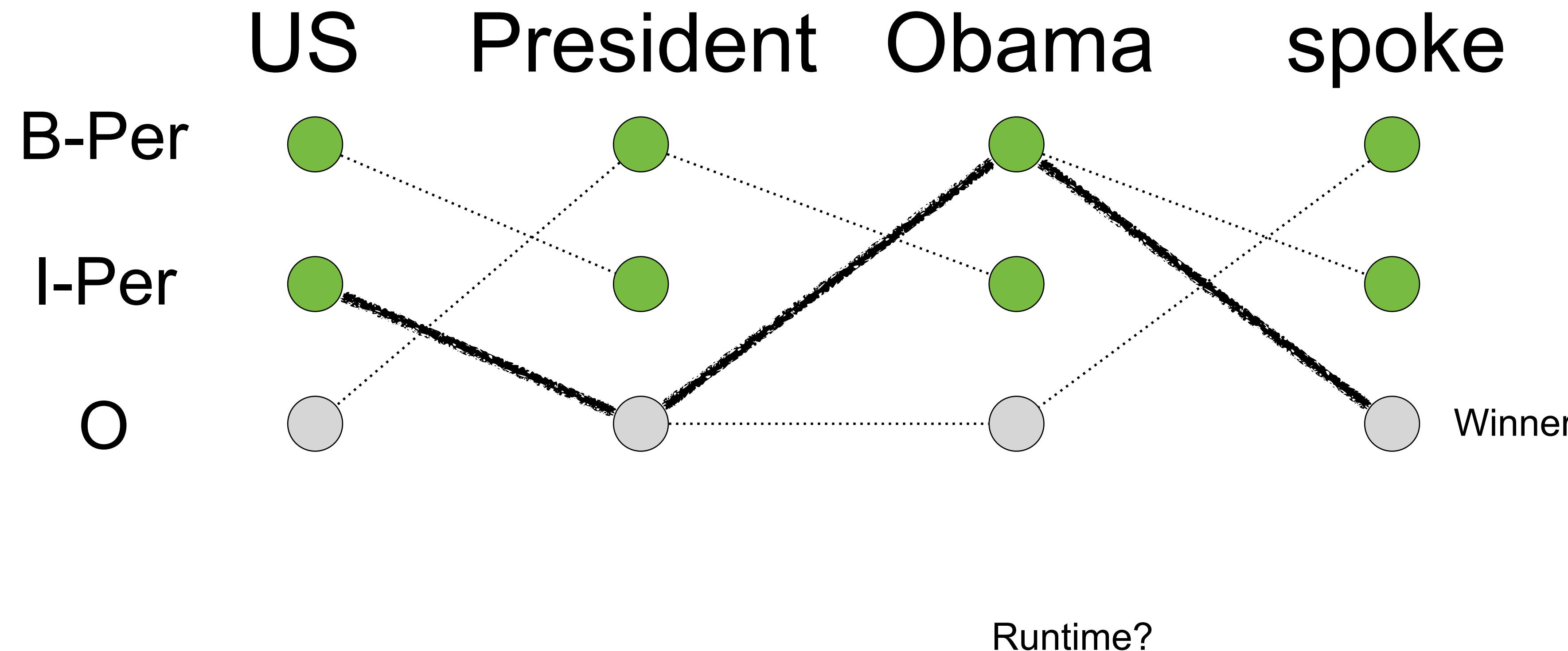
# Viterbi Decoding: Backtracking

Which was the maximizing node coming into label 3?



# Viterbi Decoding: Backtracking

Which was the maximizing node coming into label 2?



# Viterbi Summary

① Calculate scores for each token, left-to-right

- $\alpha_i(y_i) = \max_{y_{i-1}} \alpha(y_{i-1}) + (y_{i-1}, y_i)$
- Remember winning  $y_{i-1}$  for each  $i$

② Pick  $y_i$  with highest  $\alpha_i(y_i)$ , right-to-left

# Training



$$\text{LOSS} : - \sum_{\substack{(x, y) \in D}}^1 \log p(y | x)$$

We like our  
Loglikelihood loss!

we only have ↓

$$p(y_1, \dots, y_n | x) = \frac{e^{s(y_1, \dots, y_n | x)}}{\sum_{y'_1 \dots y'_n} e^{s(y'_1, \dots, y'_n | x)}} \quad \text{"Softmax" ?}$$

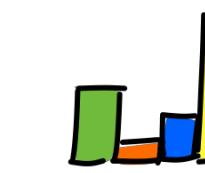
$\sum_{y'_1 \dots y'_n} e^{s(y'_1, \dots, y'_n | x)}$  =  $Z(x)$  Complexity?

Conditional Random Field

Partition Function

# Marginal Probabilities

Calculating Partition Function  $\approx$  Calculating Marg. Probabilities



Spoke at the White House

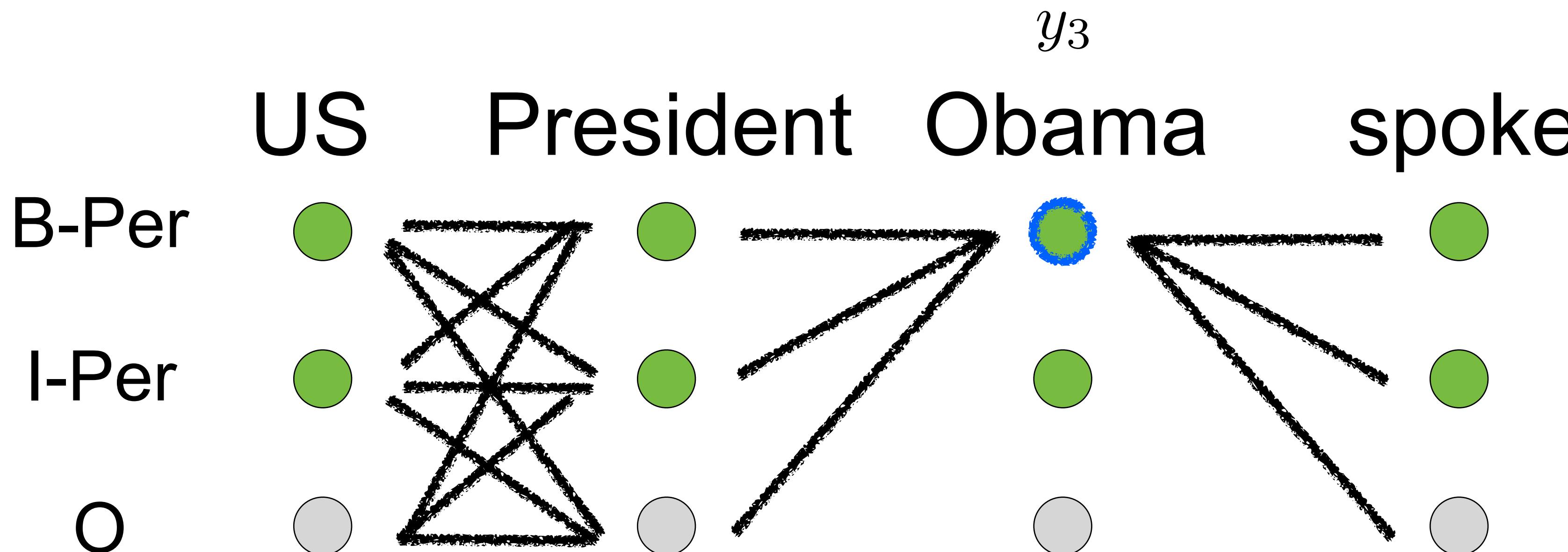
How likely is label  $y$  at position  $j$



$$P_j(y|x) = \sum_{\substack{y \\ y_j=y}} P(y|x) \propto \frac{1}{Z(x)} \sum_{\substack{y \\ y_j=y}} e^{s(y|x)}$$

# Marginal Probabilities

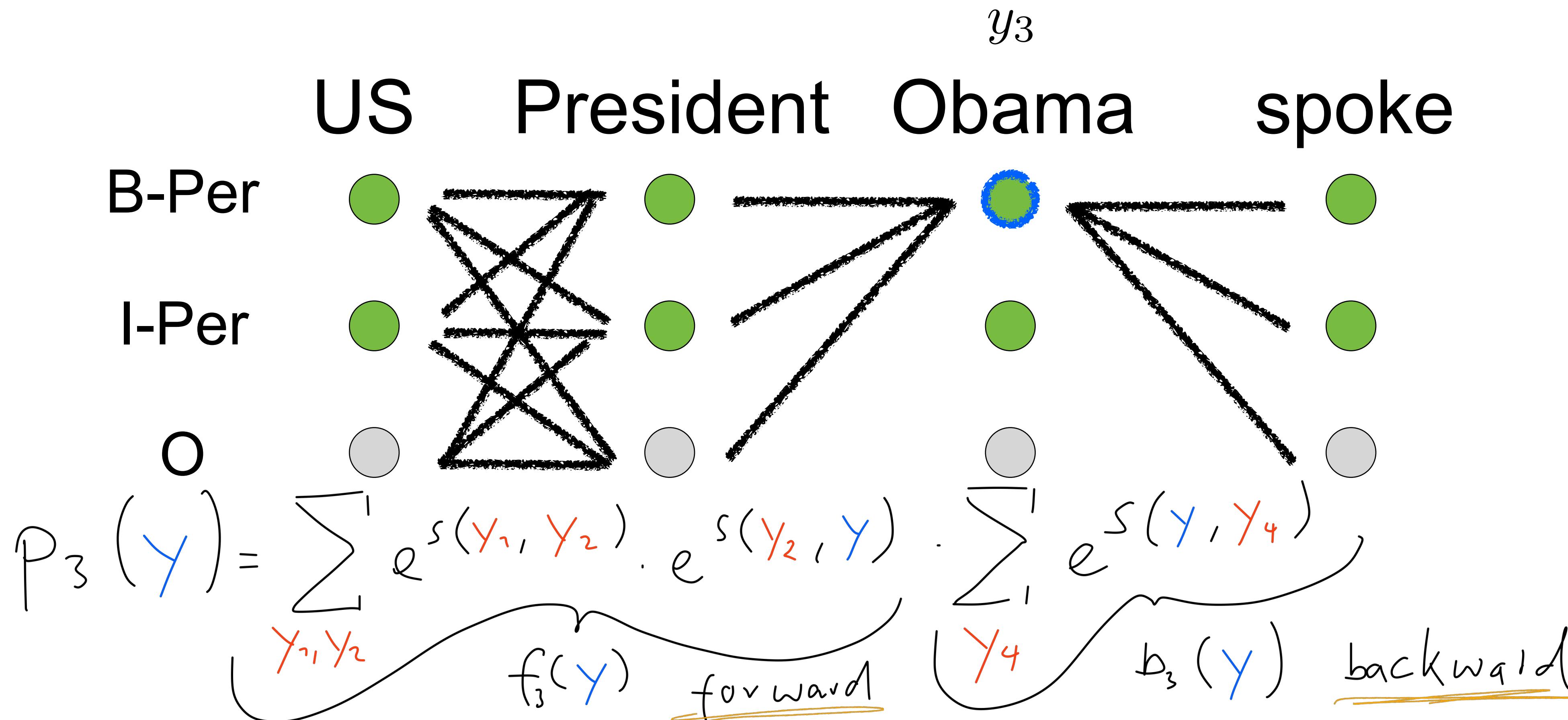
Find the marginal probability of label at token 3 to be  $y_3$



$$p_3(y) \propto \sum_{y_1, y_2, y_3=y_1, y_4} e^{s(y)} = \sum_i \prod_{i=1}^{n-1} e^{s_i(y_i, y_{i+1})}$$

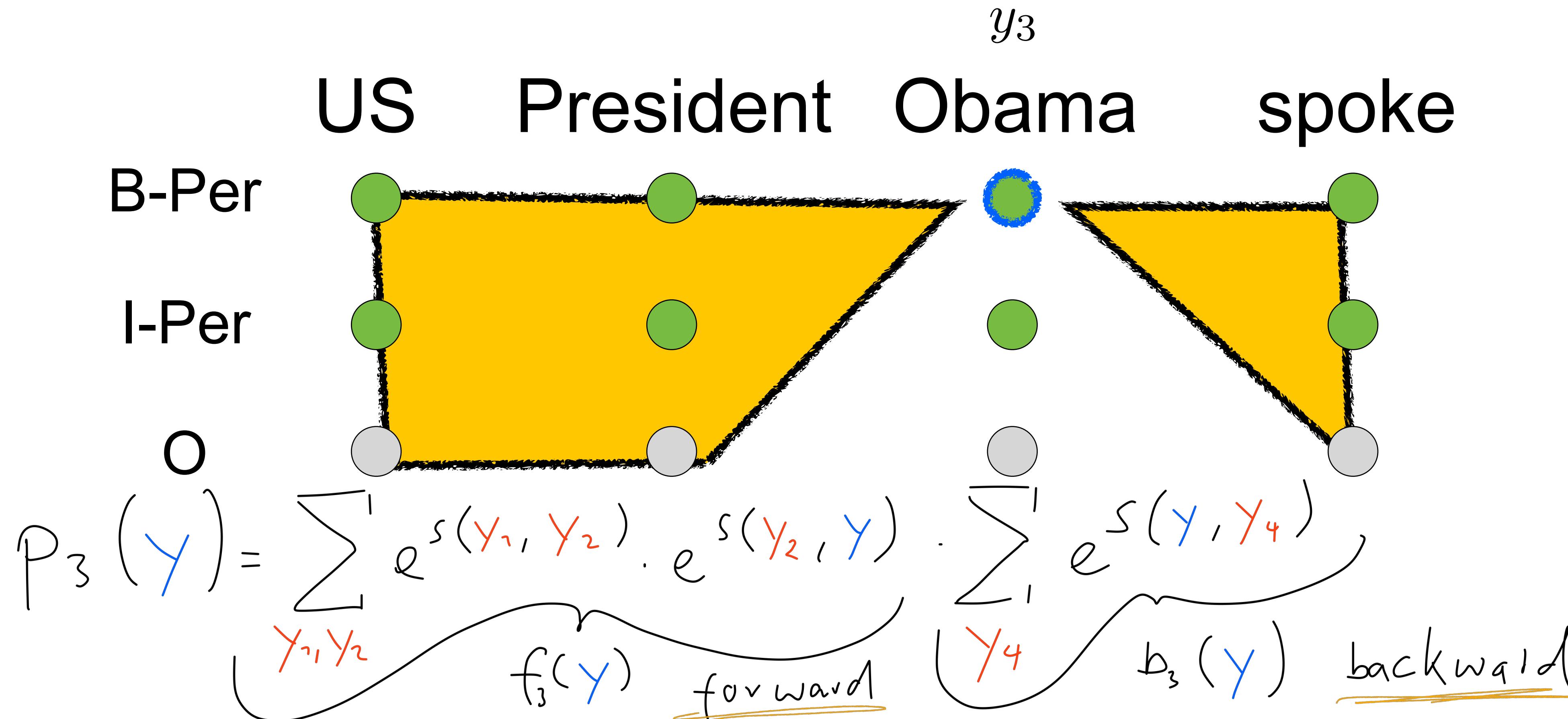
# Forward-Backward

Sum over all paths that end in  $y_3$  and start in  $y_3$



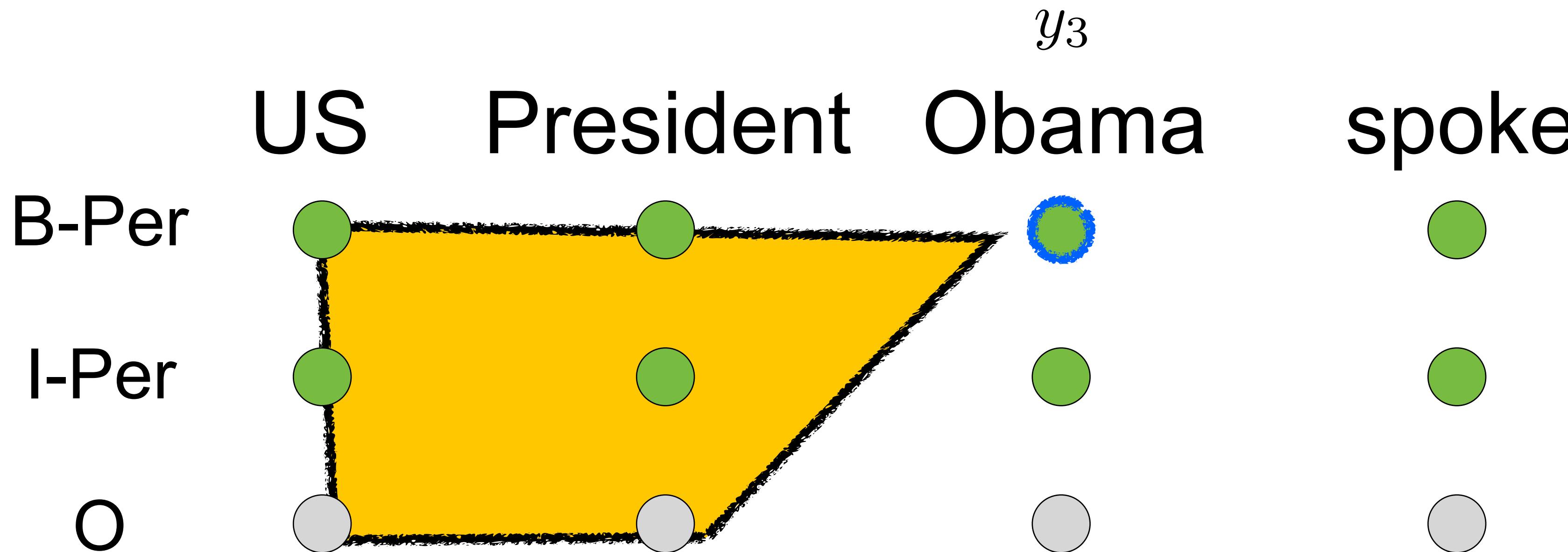
# Forward-Backward

Marginal probability of variable is product of **forward** and **backward** message



# Forward-Backward

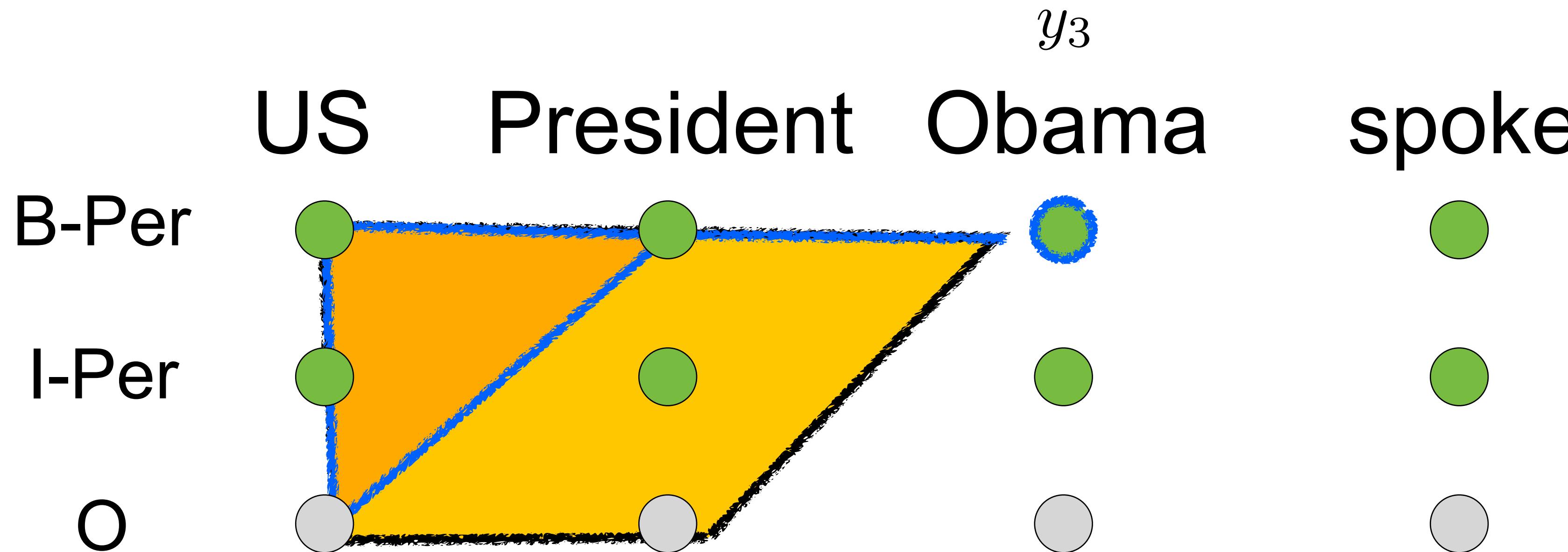
Let's focus on the forward message (backward equivalent)



$$f_j(y) = \sum_{y_1, y_2} e^{s(y_2, y)} \cdot \varrho^s(y_1, y_2)$$

# Forward-Backward

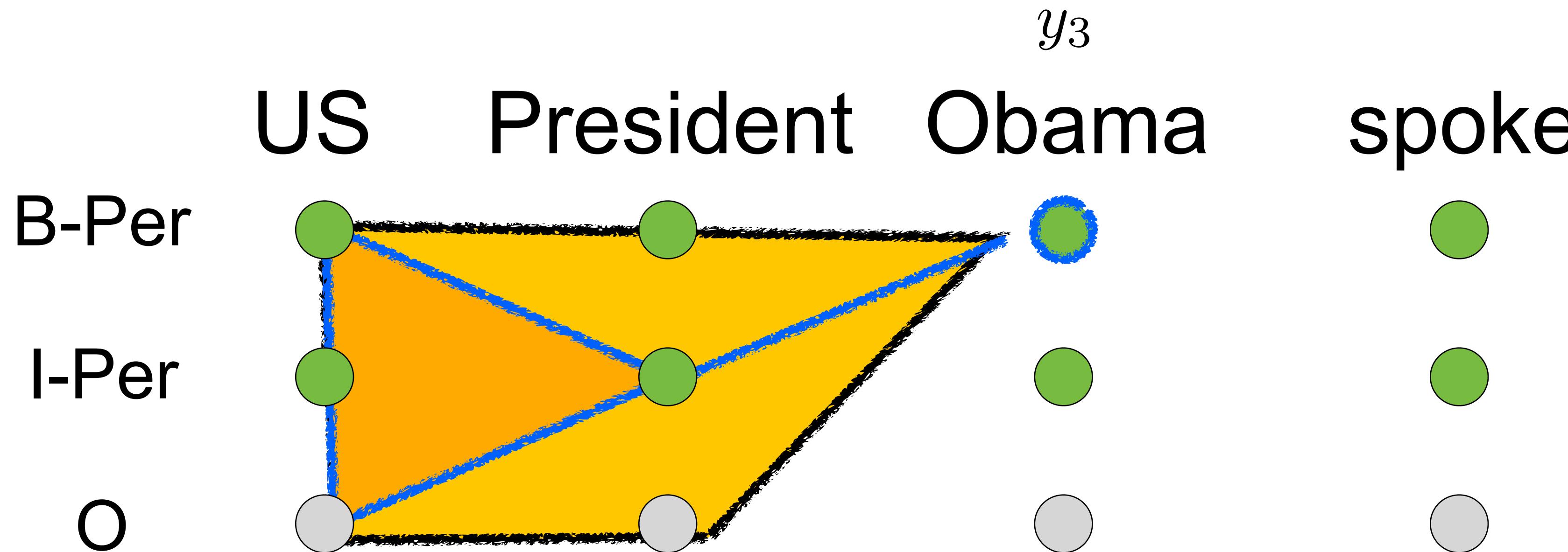
For B-PER multiply transition potential with total mass of paths going into B-PER



$$f_3(y) = \sum_{y_2} e^{s(y_2, y)} \cdot \sum_{y_1} e^{s(y_1, y_2)}$$

# Forward-Backward

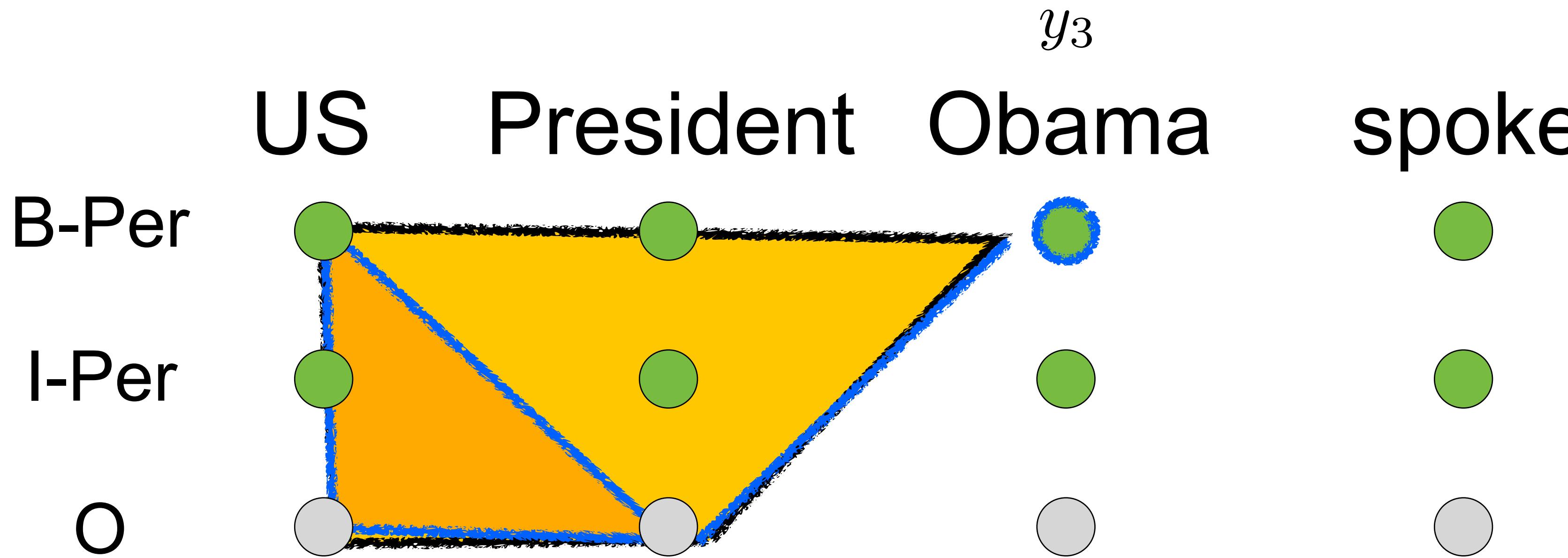
For I-PER multiply transition potential with total mass of paths going into I-PER



$$f_3(y) = \sum_{y_2} e^{s(y_2, y)} \cdot \sum_{y_1} e^{s(y_1, y_2)}$$

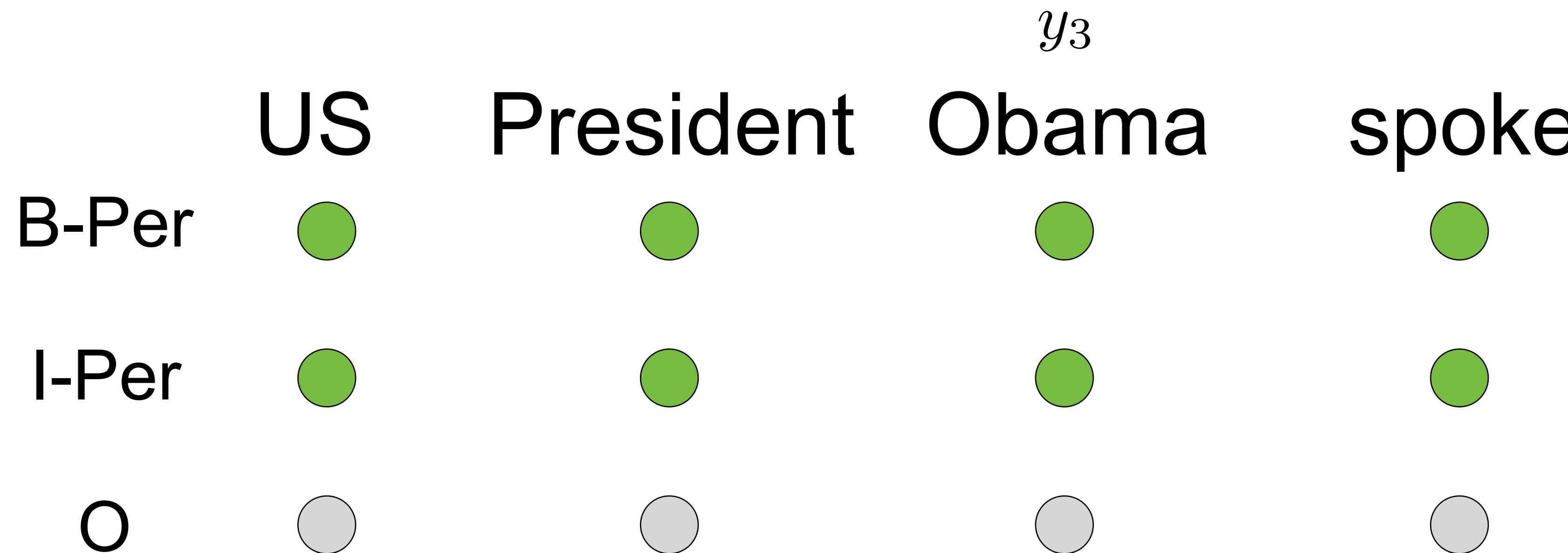
# Forward-Backward

For O multiply transition potential with total mass of paths going into O



$$f_3(y) = \sum_{y_2} e^{s(y_2, y)} \cdot \underbrace{\sum_{y_1} e^{s(y_1, y_2)}}_{f_2(y_2)}$$

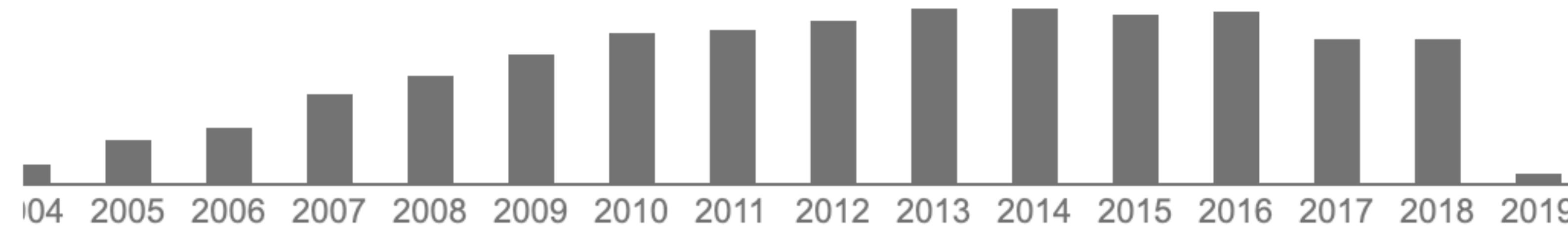
# Forward-Backward



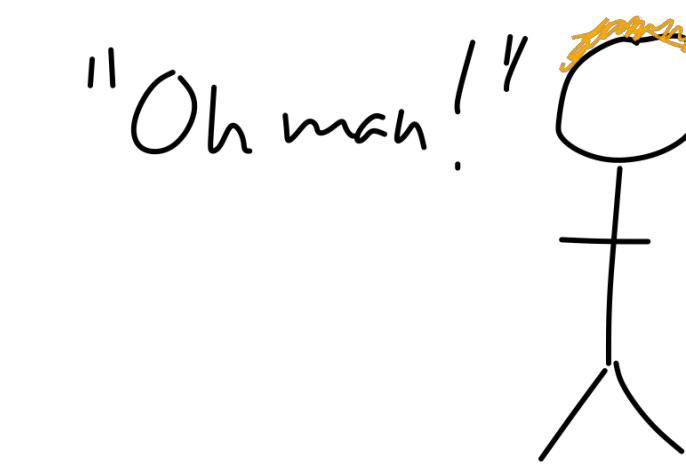
$$f_i(y) = \sum_{y_{i-1}} e^{s(y_{i-1}, y)} \cdot f_{i-1}(y_{i-1})$$

# Some History: CRF

Total citations Cited by 12054



Seb



Timster



# References

- Neural Sequence Labelling
  - J&M Chapter 9
  - Goldberg Chapter 8
  - Supervised Sequence Labelling with Recurrent Neural Networks, Alex Graves PhD Thesis
  - Design Challenges and Misconceptions in Neural Sequence Labeling, Jie Yang, Shuailong Liang, Yue Zhang, COLING 2018
- Viterbi, Forward-Backward
  - J&M Chapter 8, Appendix 3
  - Goldberg Chapter 8
  - Lafferty, J., McCallum, A., Pereira, F. (2001). "Conditional random fields: Probabilistic models for segmenting and labeling sequence data". ICML 2001