

# Information Extraction

CS6200

Information Retrieval

(and a sort of advertisement for NLP in the spring)

# Information Extraction

- Automatically extract structure from text
  - annotate document using tags to identify extracted structure
- We've briefly mentioned one example
  - But part of speech tagging is so low-level it usually doesn't count as IE
- *Named entity recognition*
  - identify words that refer to something of interest in a particular application
  - e.g., people, companies, locations, dates, product names, prices, etc.

# Named Entity Recognition

Fred Smith, who lives at 10 Water Street, Springfield, MA, is a long-time collector of **tropical fish**.

```
<p ><PersonName><GivenName>Fred</GivenName> <Sn>Smith</Sn>  
</PersonName>, who lives at <address><Street >10 Water Street</Street>,  
<City>Springfield</City>, <State>MA</State></address>, is a long-time  
collector of <b>tropical fish.</b></p>
```

- Example showing semantic annotation of text using XML tags
- Information extraction also includes document structure and more complex features such as *relationships and events*

# Named Entity Recognition

The Persian learned men say that the Phoenicians ... came to our seas from the so-called Red Sea, and having settled in the country which they still occupy, at once began to make long voyages. Among other places to which they carried Egyptian and Assyrian merchandise, they came to Argos, which was at that time preeminent in every way among the people of what is now called Hellas. The Phoenicians came to Argos, and set out their cargo. On the fifth or sixth day after their arrival, when their wares were almost all sold, many women came to the shore and among them especially the daughter of the king, whose name was Io (according to Persians and Greeks alike), the daughter of Inachus. As these stood about the stern of the ship bargaining for the wares they liked, the Phoenicians incited one another to set upon them. Most of the women escaped: Io and others were seized and thrown into the ship, which then sailed away for Egypt.

# Named Entity Recognition

The Persian learned men say that the Phoenicians ... came to our seas from the so-called Red Sea, and having settled in the country which they still occupy, at once began to make long voyages. Among other places to which they carried Egyptian and Assyrian merchandise, they came to Argos, which was at that time preeminent in every way among the people of what is now called Hellas. The Phoenicians came to Argos, and set out their cargo. On the fifth or sixth day after their arrival, when their wares were almost all sold, many women came to the shore and among them especially the daughter of the king, whose name was **Io** (according to Persians and Greeks alike), the daughter of **Inachus**. As these stood about the stern of the ship bargaining for the wares they liked, the Phoenicians incited one another to set upon them. Most of the women escaped: **Io** and others were seized and thrown into the ship, which then sailed away for Egypt.

**Person**

# Named Entity Recognition

The Persian learned men say that the Phoenicians ... came to our seas from the so-called **Red Sea**, and having settled in the country which they still occupy, at once began to make long voyages. Among other places to which they carried Egyptian and Assyrian merchandise, they came to **Argos**, which was at that time preeminent in every way among the people of what is now called **Hellas**. The Phoenicians came to **Argos**, and set out their cargo. On the fifth or sixth day after their arrival, when their wares were almost all sold, many women came to the shore and among them especially the daughter of the king, whose name was **Io** (according to Persians and Greeks alike), the daughter of **Inachus**. As these stood about the stern of the ship bargaining for the wares they liked, the Phoenicians incited one another to set upon them. Most of the women escaped: **Io** and others were seized and thrown into the ship, which then sailed away for **Egypt**.

**Person**

**Location**

# Named Entity Recognition

The **Persian** learned men say that the **Phoenicians** ... came to our seas from the so-called **Red Sea**, and having settled in the country which they still occupy, at once began to make long voyages. Among other places to which they carried **Egyptian** and **Assyrian** merchandise, they came to **Argos**, which was at that time preeminent in every way among the people of what is now called **Hellas**. The **Phoenicians** came to **Argos**, and set out their cargo. On the fifth or sixth day after their arrival, when their wares were almost all sold, many women came to the shore and among them especially the daughter of the king, whose name was **Io** (according to **Persians** and **Greeks** alike), the daughter of **Inachus**. As these stood about the stern of the ship bargaining for the wares they liked, the **Phoenicians** incited one another to set upon them. Most of the women escaped: **Io** and others were seized and thrown into the ship, which then sailed away for **Egypt**.

**Person**

**Location**

**Ethnic**

# Named Entity Recognition

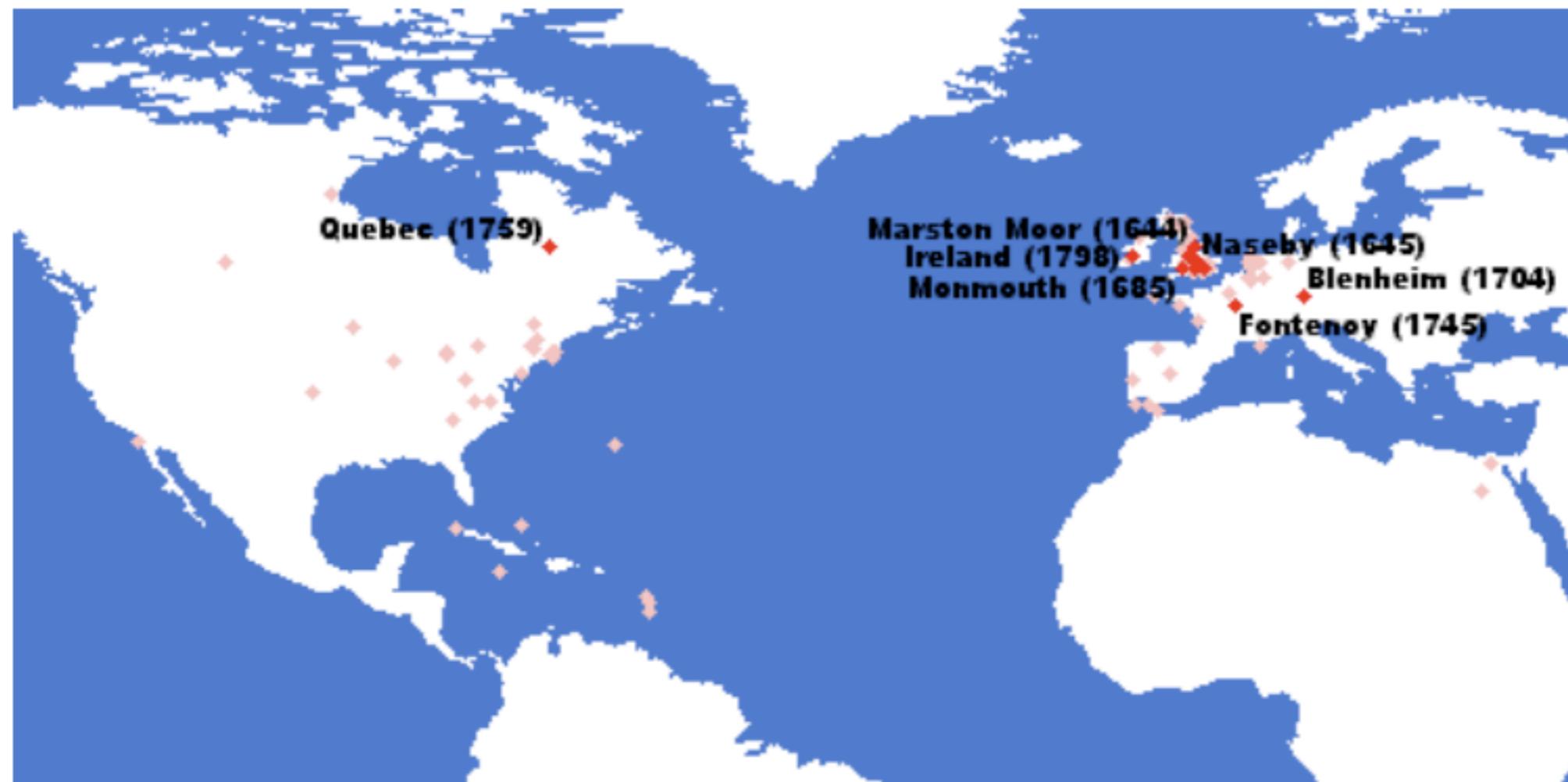
The **Persian** learned men say that the **Phoenicians** ... came to our seas from the so-called **Red Sea**, and having settled in the country which they still occupy, at once began to make long voyages. Among other places to which they carried **Egyptian** and **Assyrian** merchandise, they came to **Argos**, which was at that time preeminent in every way among the people of what is now called **Hellas**. The **Phoenicians** came to **Argos**, and set out their cargo. On the fifth or sixth day after their arrival, when their wares were almost all sold, many women came to the shore and among them especially the daughter of the king, whose name was **Io** (according to **Persians** and **Greeks** alike), the daughter of **Inachus**. As these stood about the stern of the ship bargaining for the wares they liked, the **Phoenicians** incited one another to set upon them. Most of the women escaped: **Io** and others were seized and thrown into the ship, which then sailed away for **Egypt**.

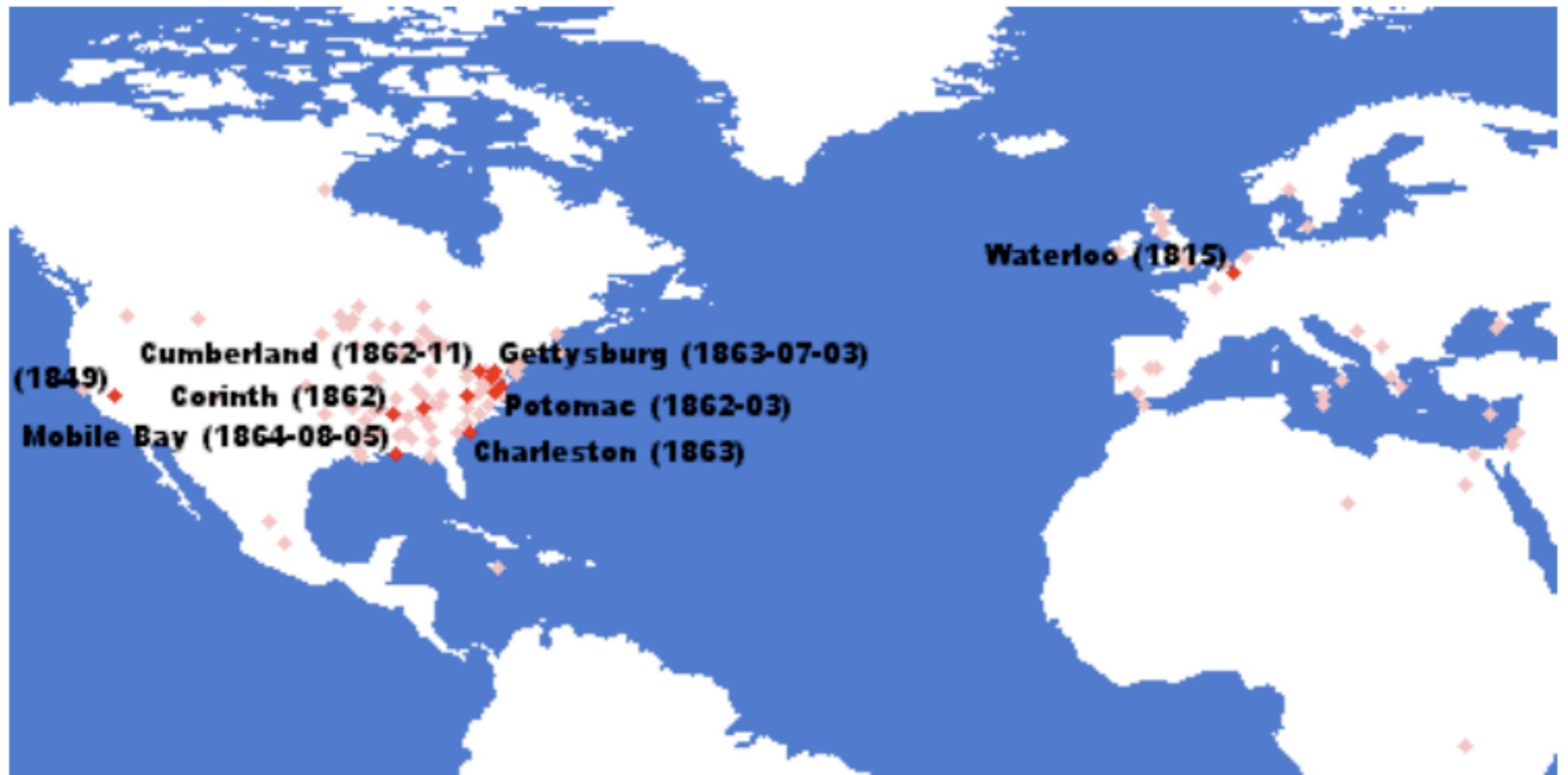
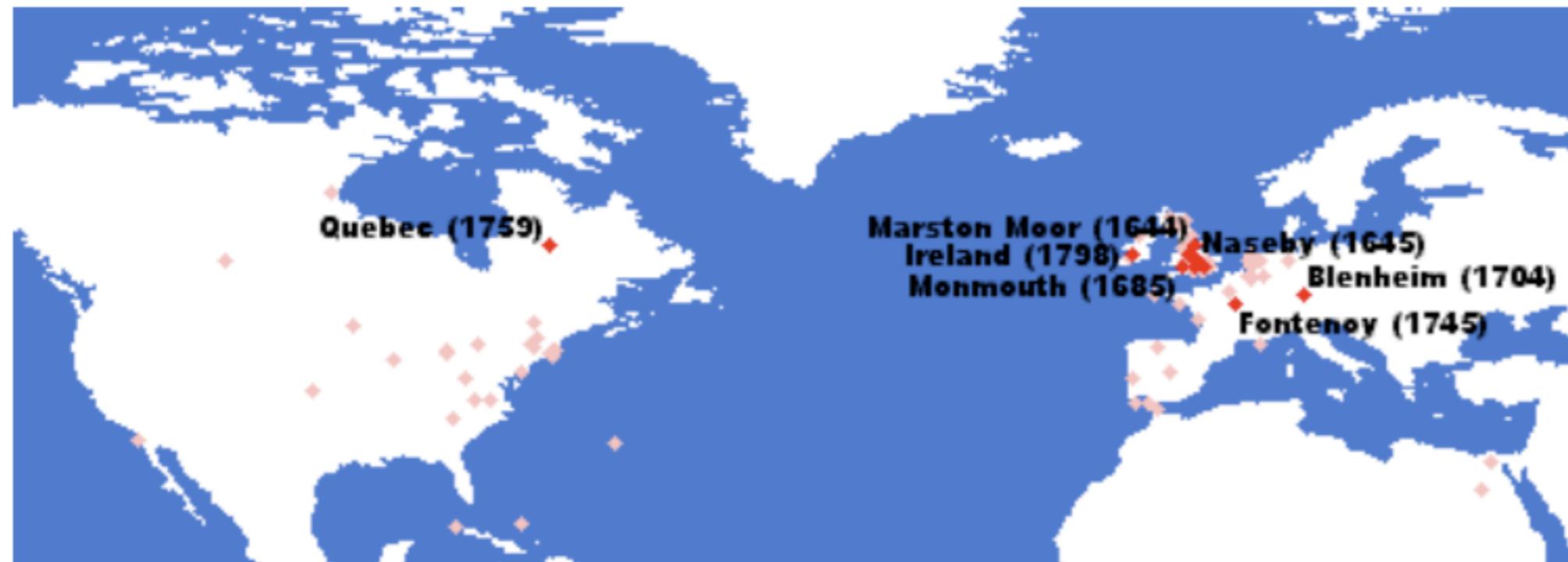
*Classes could also be, e.g., Wikipedia articles*

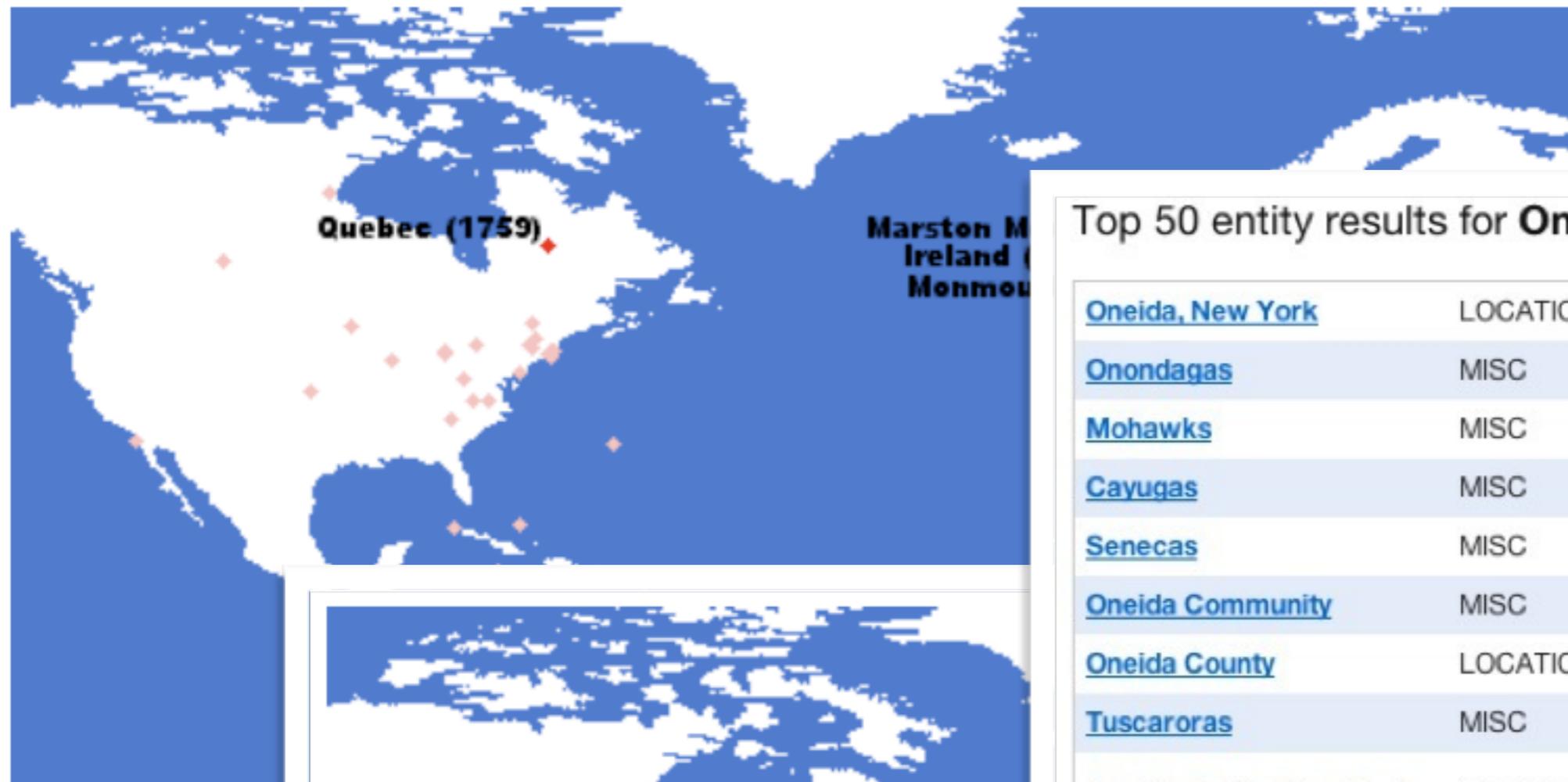
**Person**

**Location**

**Ethnic**







### Top 50 entity results for Oneida

<a href="#">Oneida, New York</a>	LOCATION	(Longitude: -75
<a href="#">Onondagas</a>	MISC	
<a href="#">Mohawks</a>	MISC	
<a href="#">Cayugas</a>	MISC	
<a href="#">Senecas</a>	MISC	
<a href="#">Oneida Community</a>	MISC	
<a href="#">Oneida County</a>	LOCATION	(Longitude: -75
<a href="#">Tuscaroras</a>	MISC	
<a href="#">Oneida Castle, New York</a>	LOCATION	(Longitude: -75 43.078333333333
<a href="#">Oneida Conference</a>	ORGANIZATION	
<a href="#">Oneida Indians</a>	MISC	
<a href="#">Mohicans</a>	MISC	
<a href="#">Mohawk</a>	MISC	

# Named Entity Recognition

- *Rule-based*
  - Uses *lexicons* (lists of words and phrases) that categorize names
    - e.g., locations, peoples' names, organizations, etc.
  - Rules also used to verify or find new entity names
    - e.g., “<number> <word> street” for addresses
    - “<street address>, <city>” or “in <city>” to verify city names
    - “<street address>, <city>, <state>” to find new cities
    - “<title> <name>” to find new names

# Named Entity Recognition

- Rules either developed manually by trial and error or using machine learning techniques
- *Statistical*
  - uses a probabilistic model of the words in and around an entity
  - probabilities estimated using *training data* (manually annotated text)
  - Hidden Markov Model (HMM) is one approach
    - Conditional Random Fields: similar structure, often higher accuracy, more expensive to train

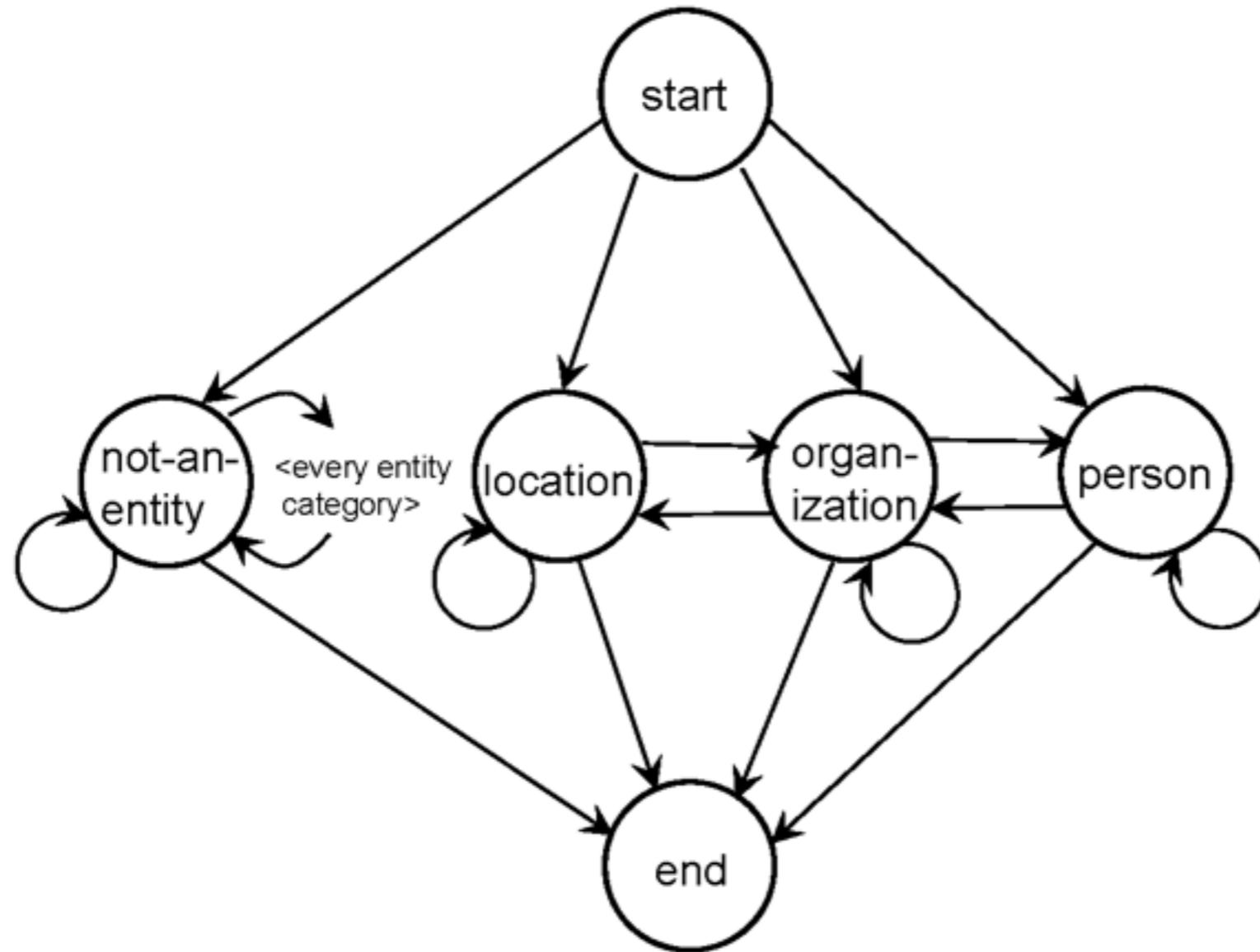
# HMM for Extraction

- Resolve ambiguity in a word using *context*
  - e.g., “marathon” is a location or a sporting event, “boston marathon” is a specific sporting event
- Model context using a *generative* model of the sequence of words
  - *Markov property*: the next word in a sequence depends only on a small number of the previous words

# HMM for Extraction

- *Markov Model* describes a process as a collection of states with transitions between them
  - each transition has a probability associated with it
  - next state depends only on current state and transition probabilities
- *Hidden Markov Model*
  - each state has a set of possible outputs
  - outputs have probabilities

# HMM Sentence Model



- Each state is associated with a probability distribution over words (the output)

# NER as Sequence Tagging

The Phoenicians came from the Red Sea

# NER as Sequence Tagging

○      **B-E**      ○      ○      ○      **B-L**      **I-L**  
The    Phoenicians    came    from    the    Red    Sea

# Sequence Tagging

Fed raises interest rates

# Sequence Tagging

NN

NNS

NNP

VB

VBZ

Fed

raises

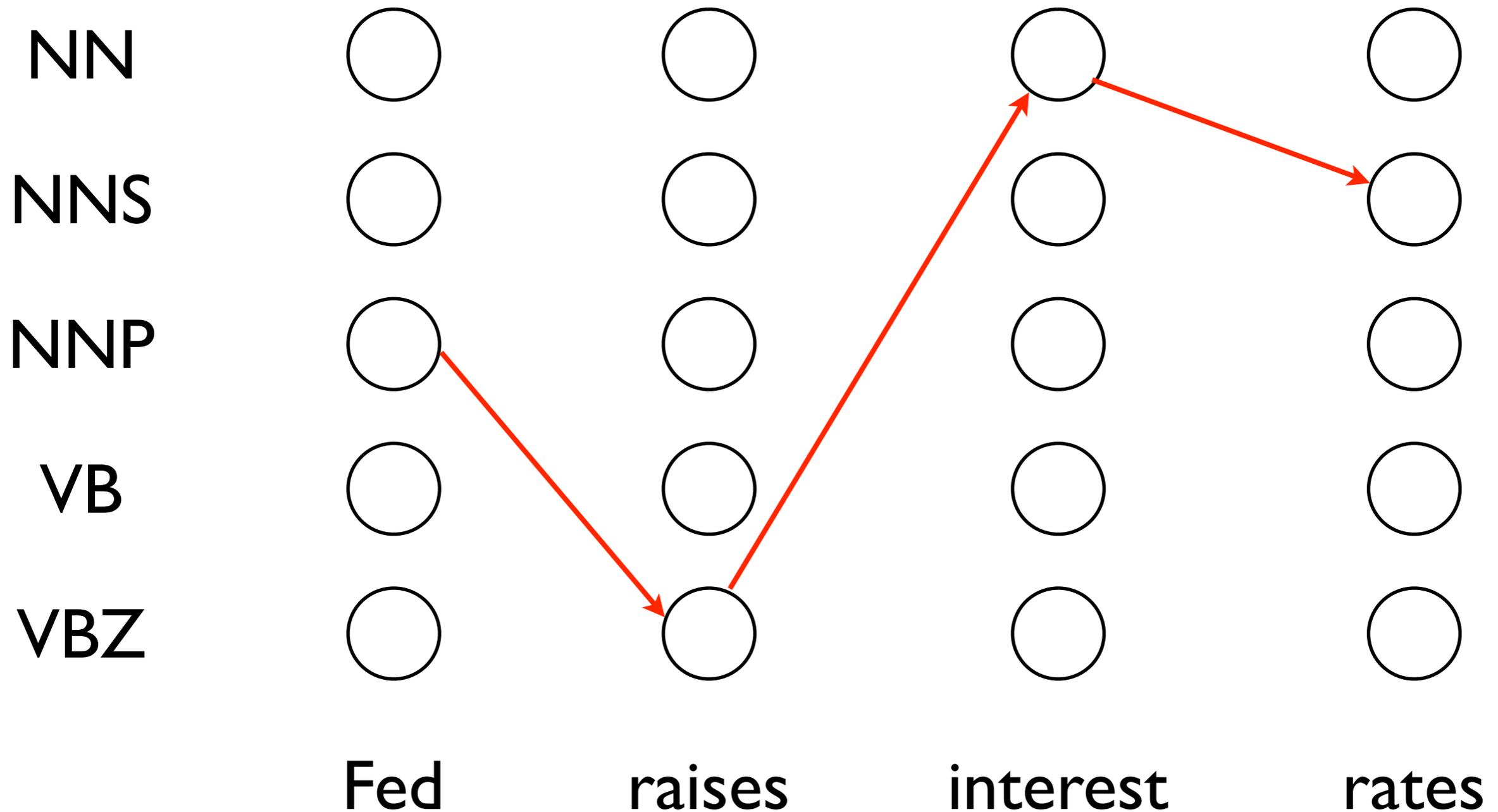
interest

rates

# Sequence Tagging

NN	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
NNS	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
NNP	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
VB	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
VBZ	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
	Fed	raises	interest	rates

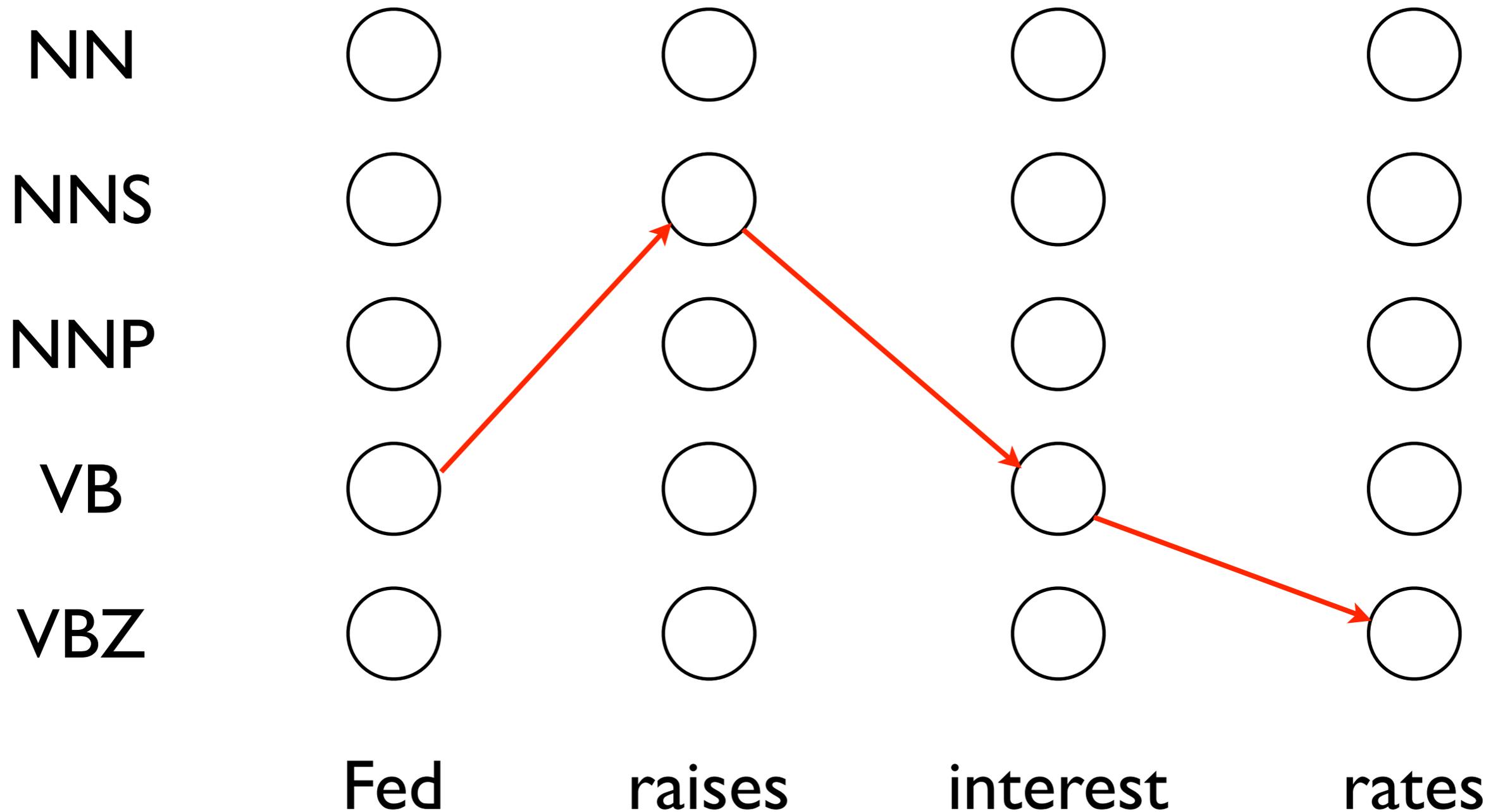
# Sequence Tagging



# Sequence Tagging

NN	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
NNS	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
NNP	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
VB	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
VBZ	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
	Fed	raises	interest	rates

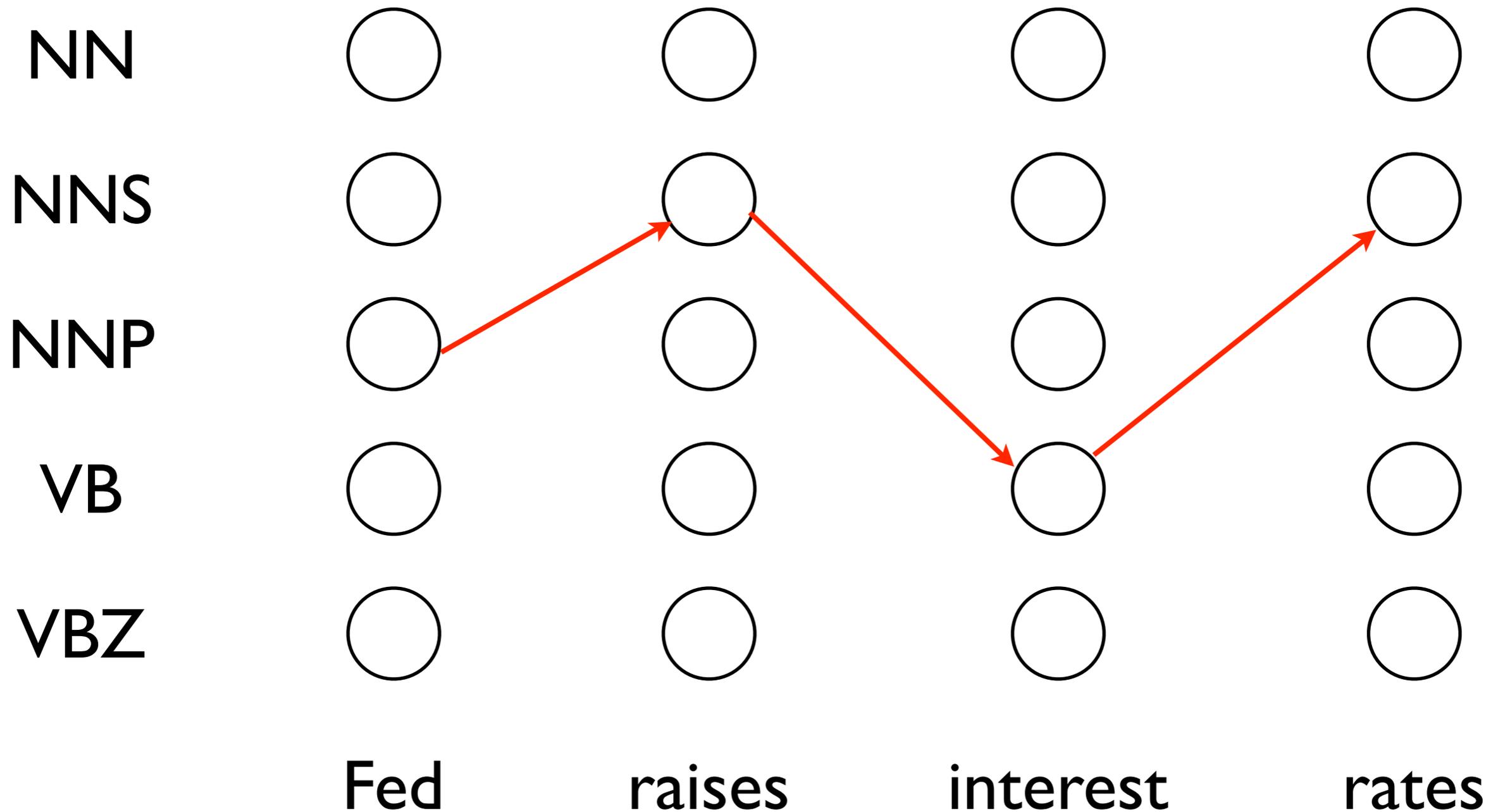
# Sequence Tagging



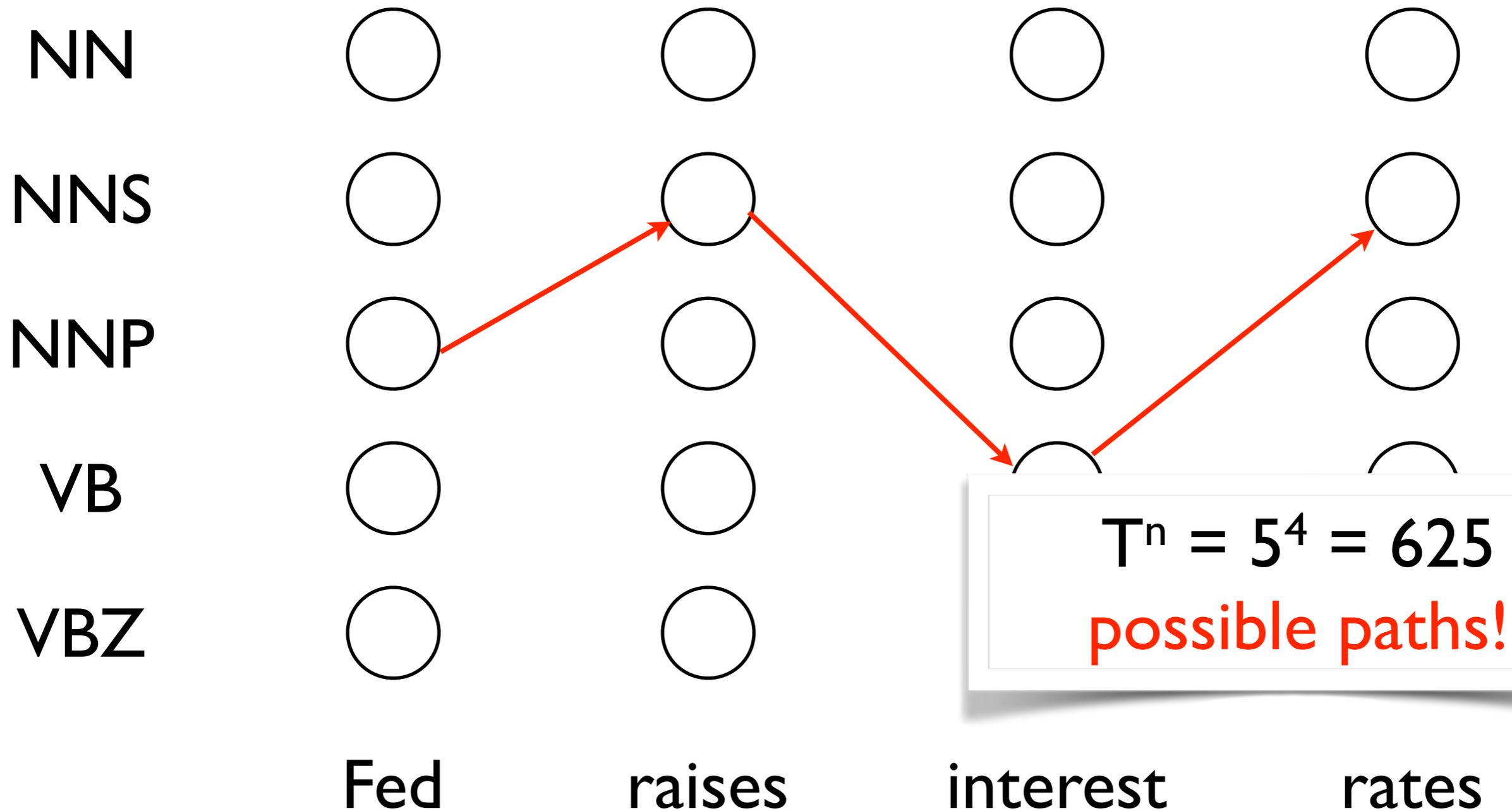
# Sequence Tagging

NN	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
NNS	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
NNP	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
VB	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
VBZ	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
	Fed	raises	interest	rates

# Sequence Tagging

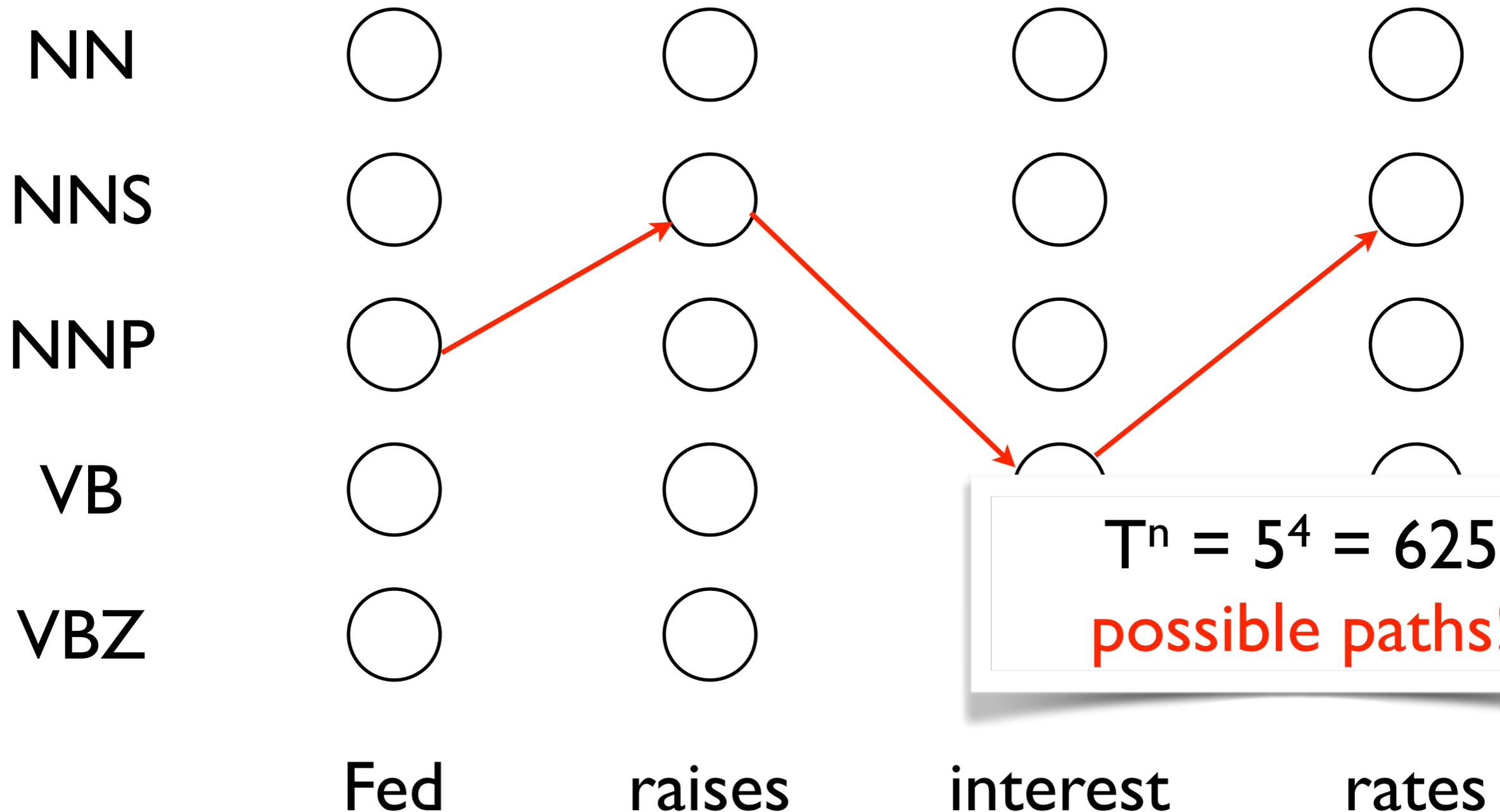


# Sequence Tagging



# Sequence Tagging

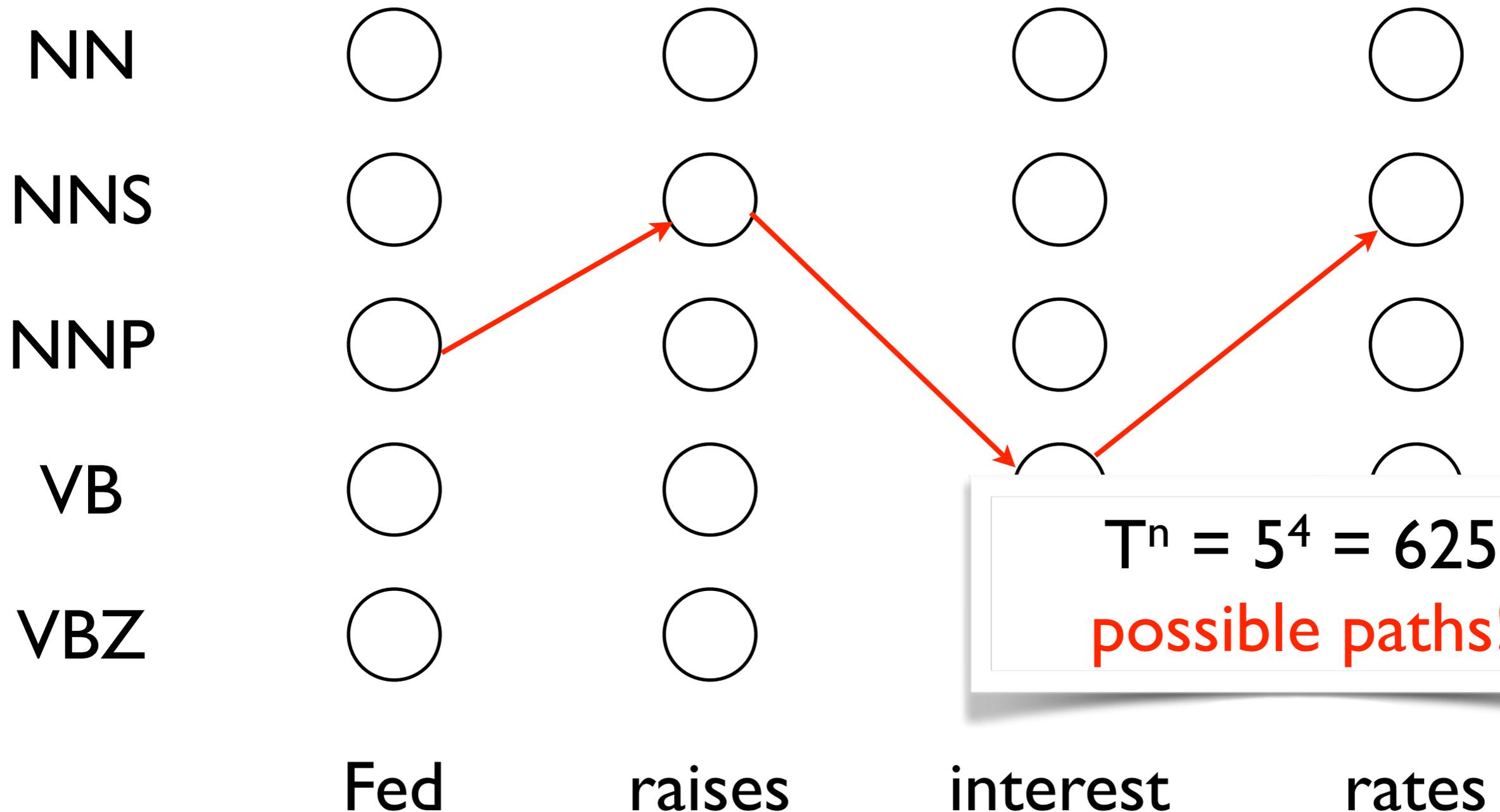
*Efficient (linear time) Shortest path = Viterbi algorithm*



$T^n = 5^4 = 625$   
possible paths!

# Sequence Tagging

*Efficient (linear time) Shortest path = Viterbi algorithm*



*Can we specify that "Fed" always has the same tag in this document?*

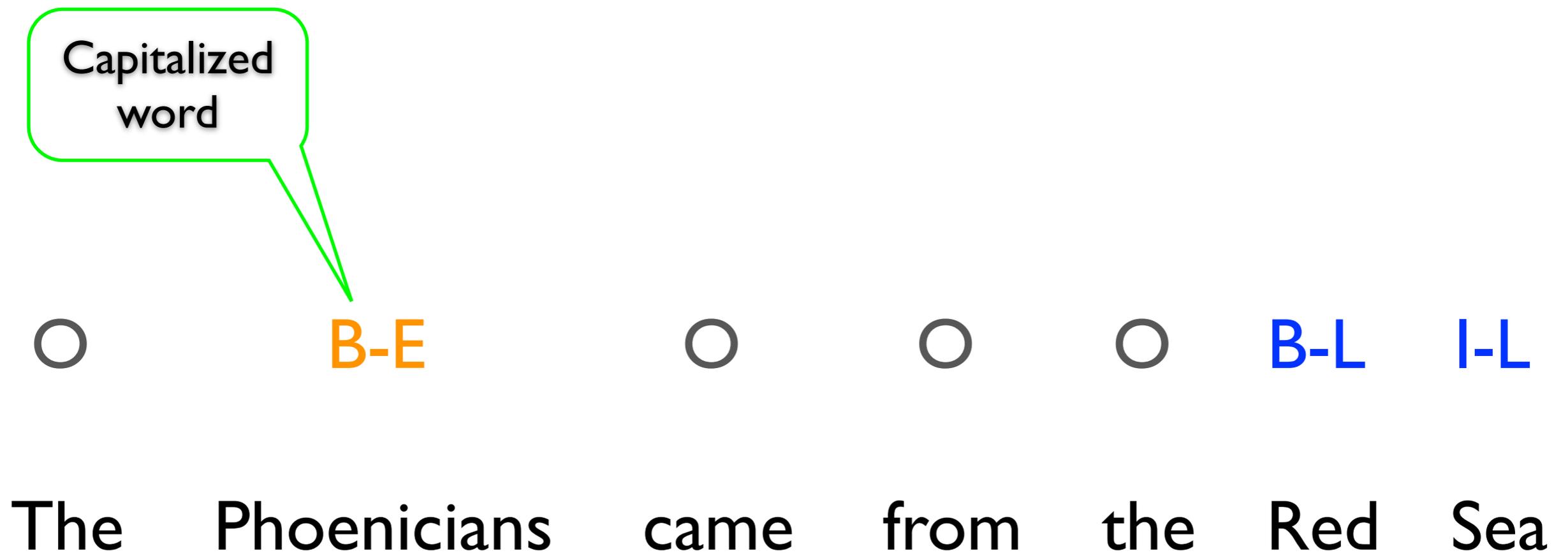
# NER as Sequence Tagging

The Phoenicians came from the Red Sea

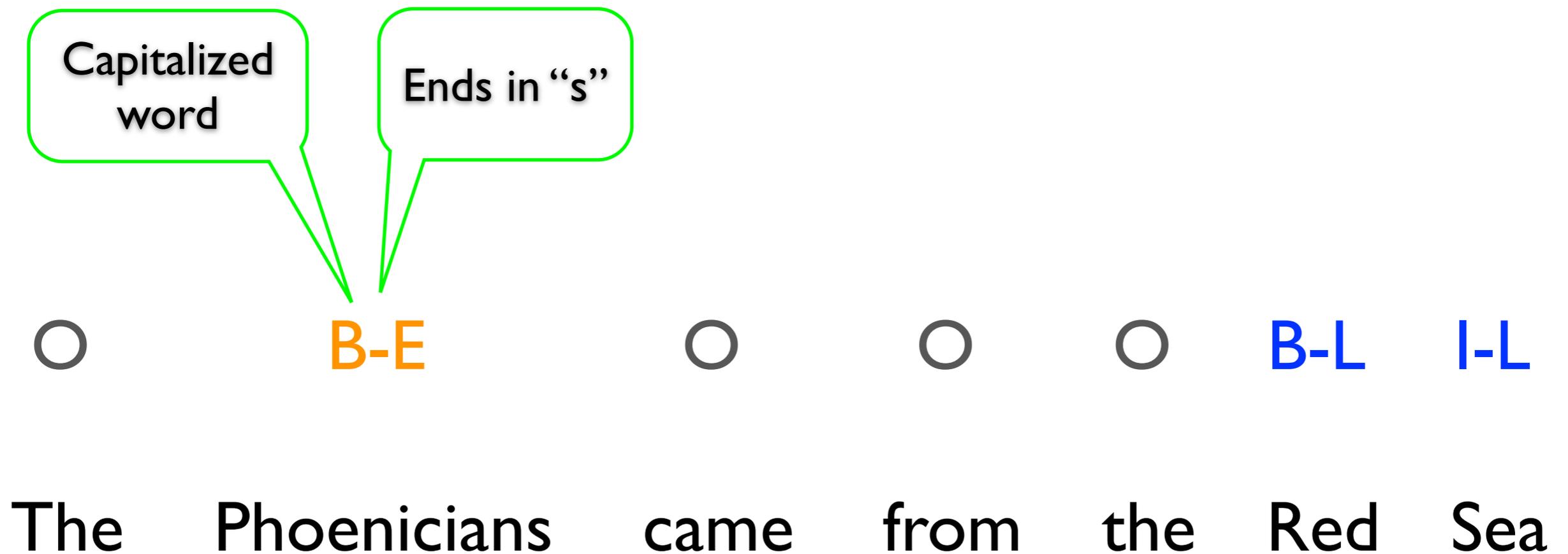
# NER as Sequence Tagging

○      **B-E**      ○      ○      ○      **B-L**      **I-L**  
The    Phoenicians    came    from    the    Red    Sea

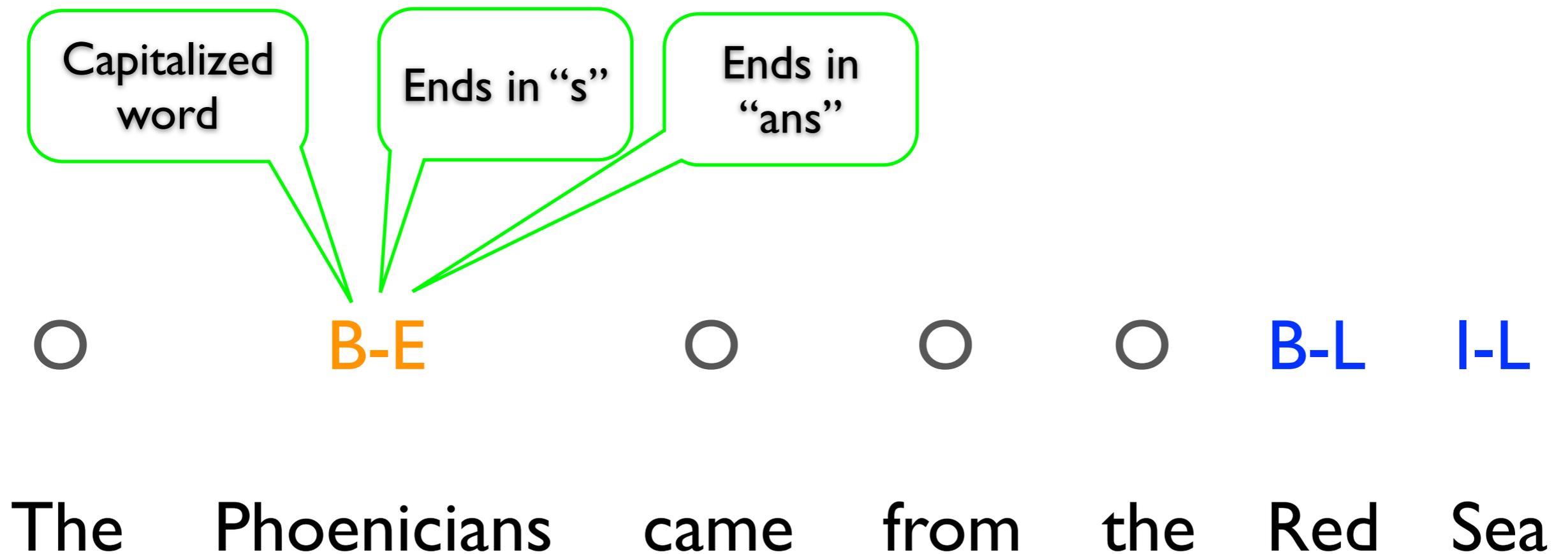
# NER as Sequence Tagging



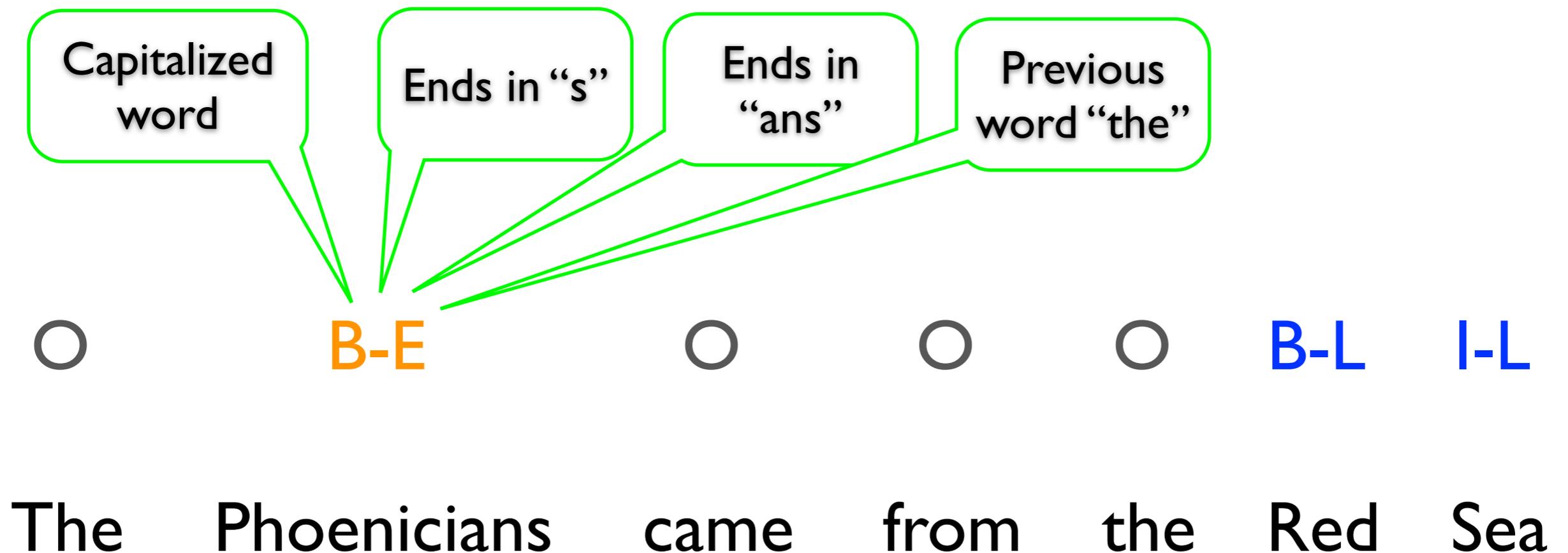
# NER as Sequence Tagging



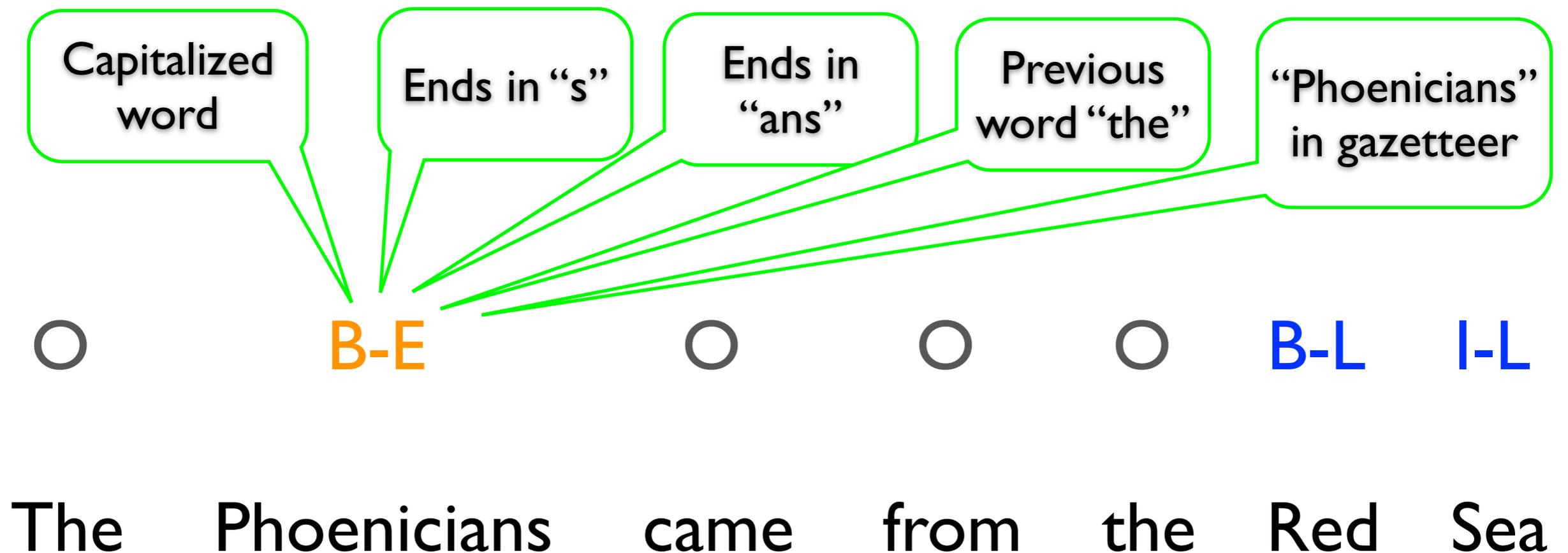
# NER as Sequence Tagging



# NER as Sequence Tagging



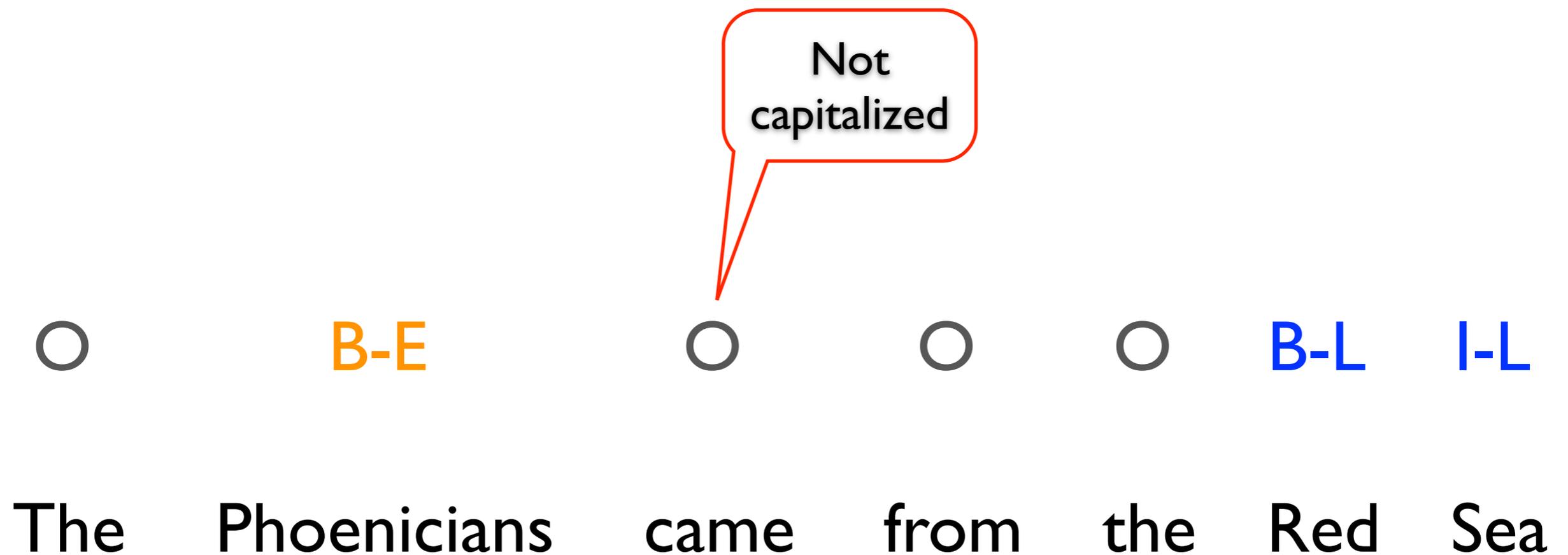
# NER as Sequence Tagging



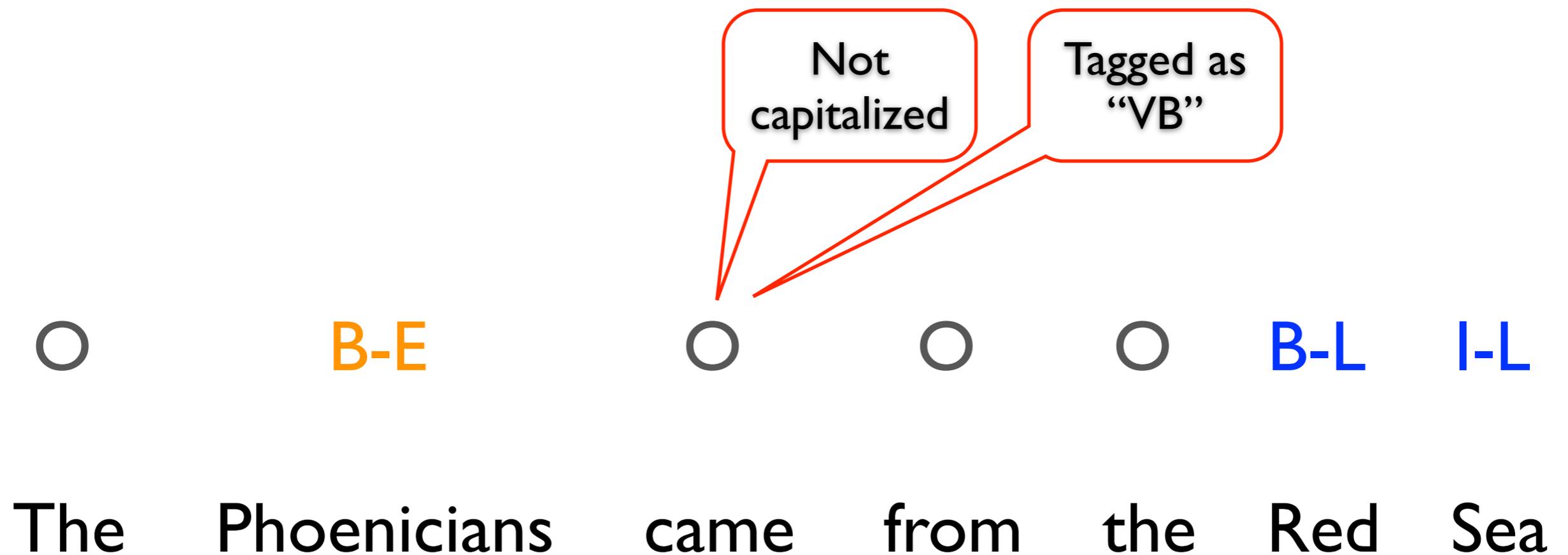
# NER as Sequence Tagging

○      **B-E**      ○      ○      ○      **B-L**      **I-L**  
The    Phoenicians    came    from    the    Red    Sea

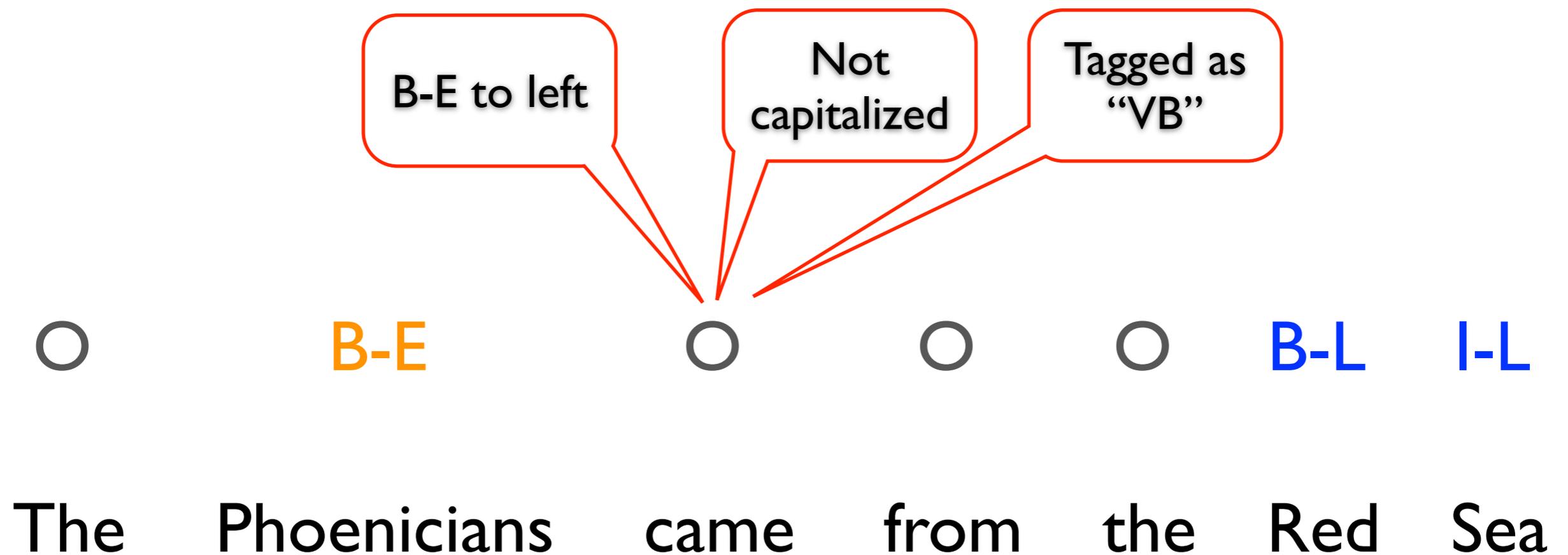
# NER as Sequence Tagging



# NER as Sequence Tagging



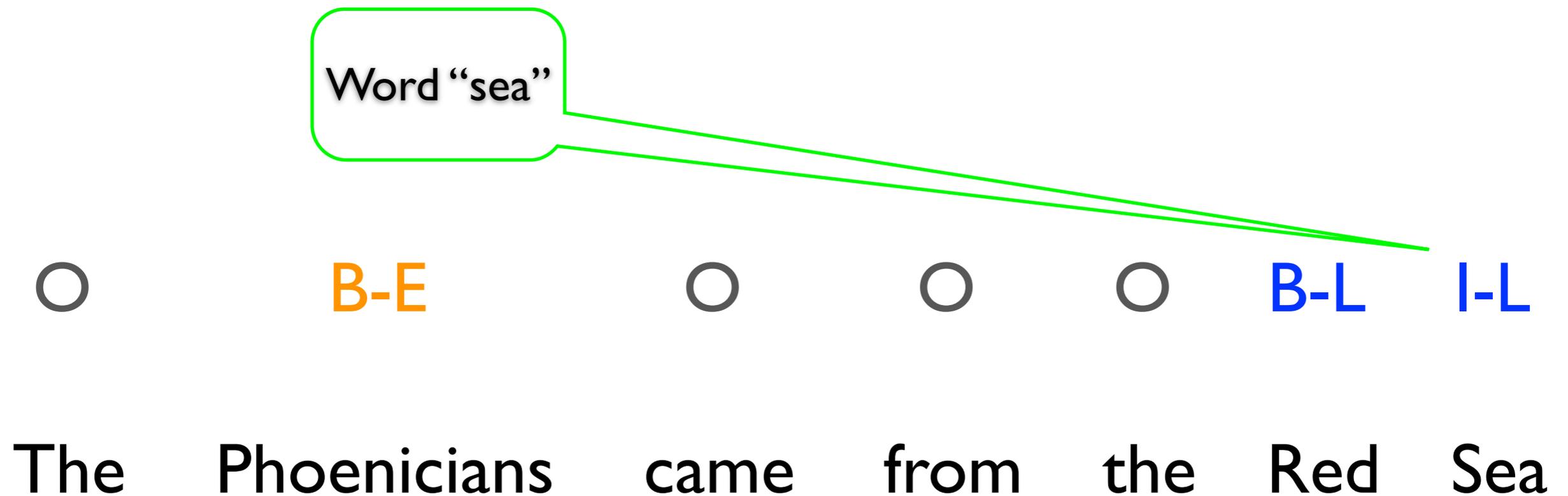
# NER as Sequence Tagging



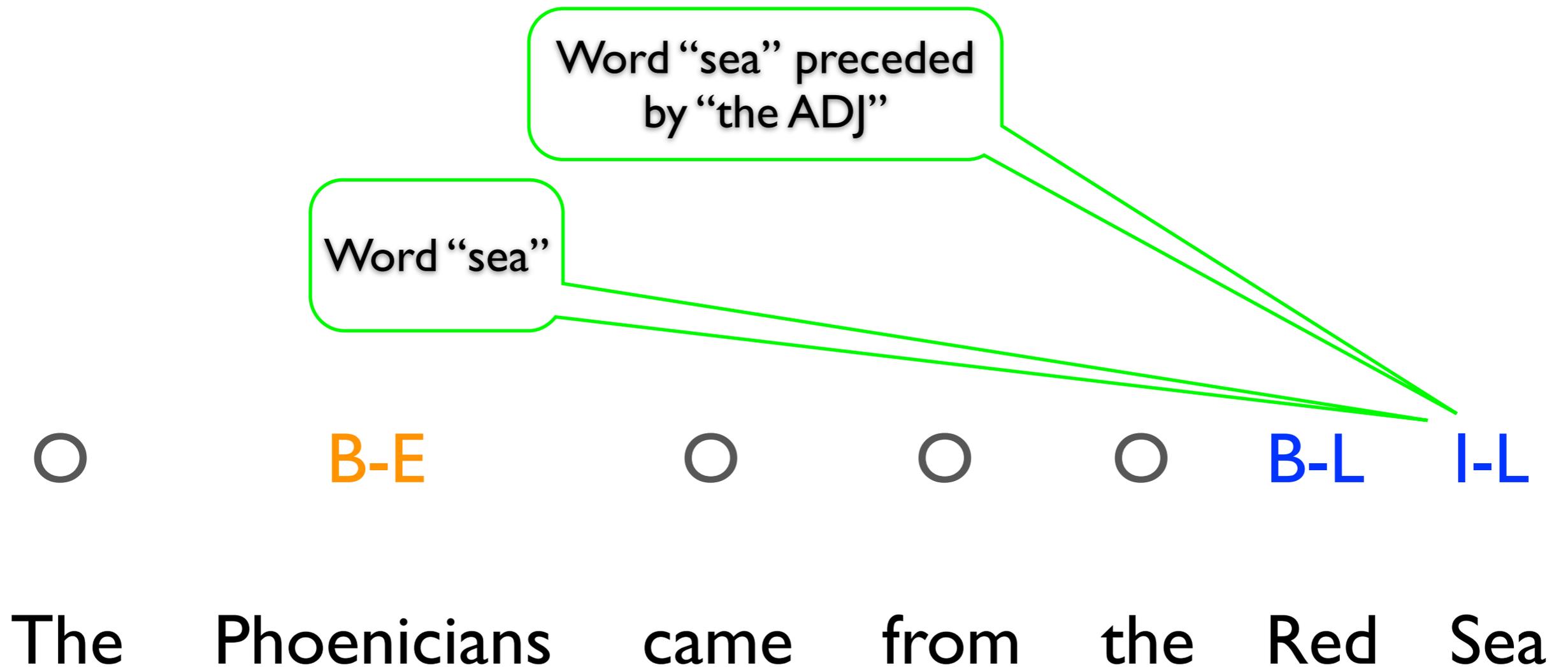
# NER as Sequence Tagging

○      **B-E**      ○      ○      ○      **B-L**      **I-L**  
The    Phoenicians    came    from    the    Red    Sea

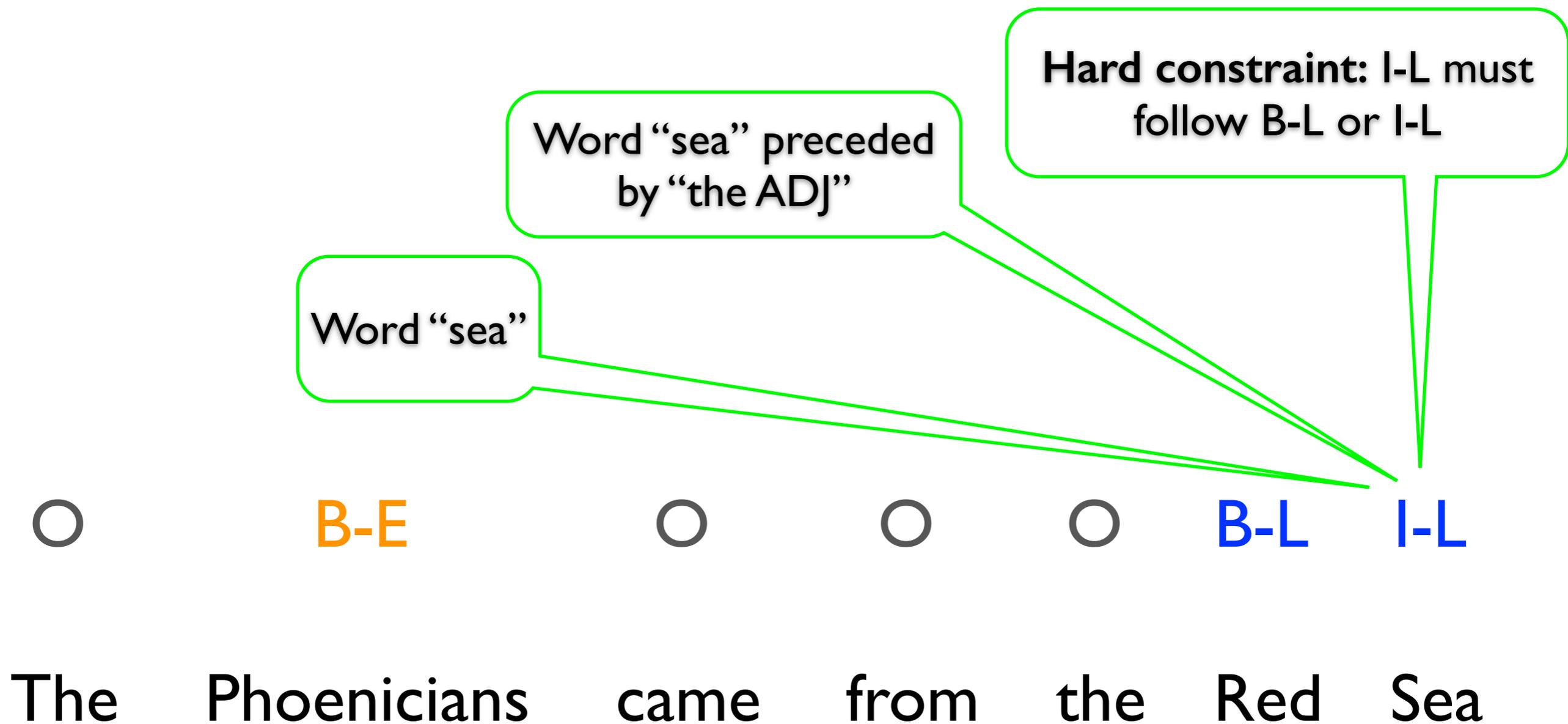
# NER as Sequence Tagging



# NER as Sequence Tagging



# NER as Sequence Tagging



# Great Ideas in ML: Message Passing



# Great Ideas in ML: Message Passing

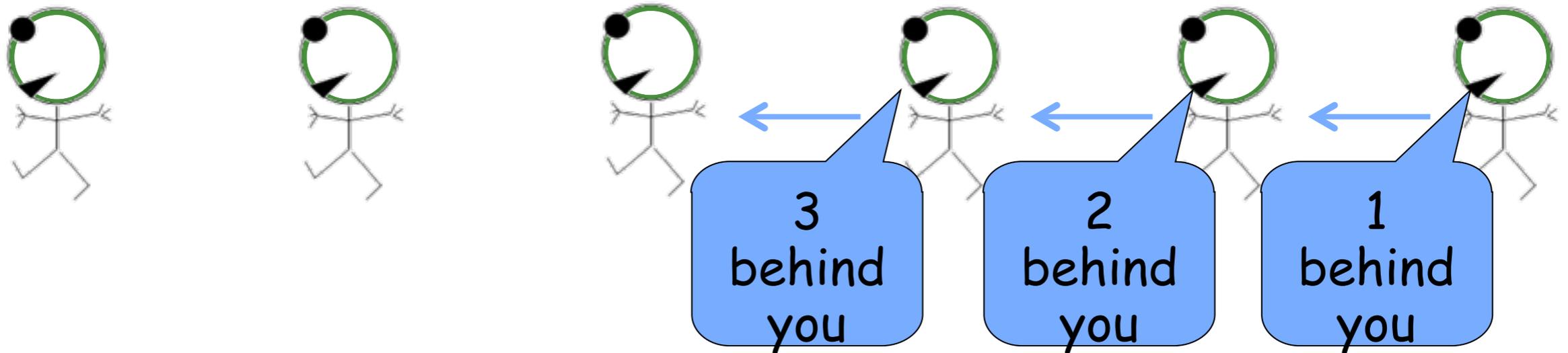
*Count the soldiers*



adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

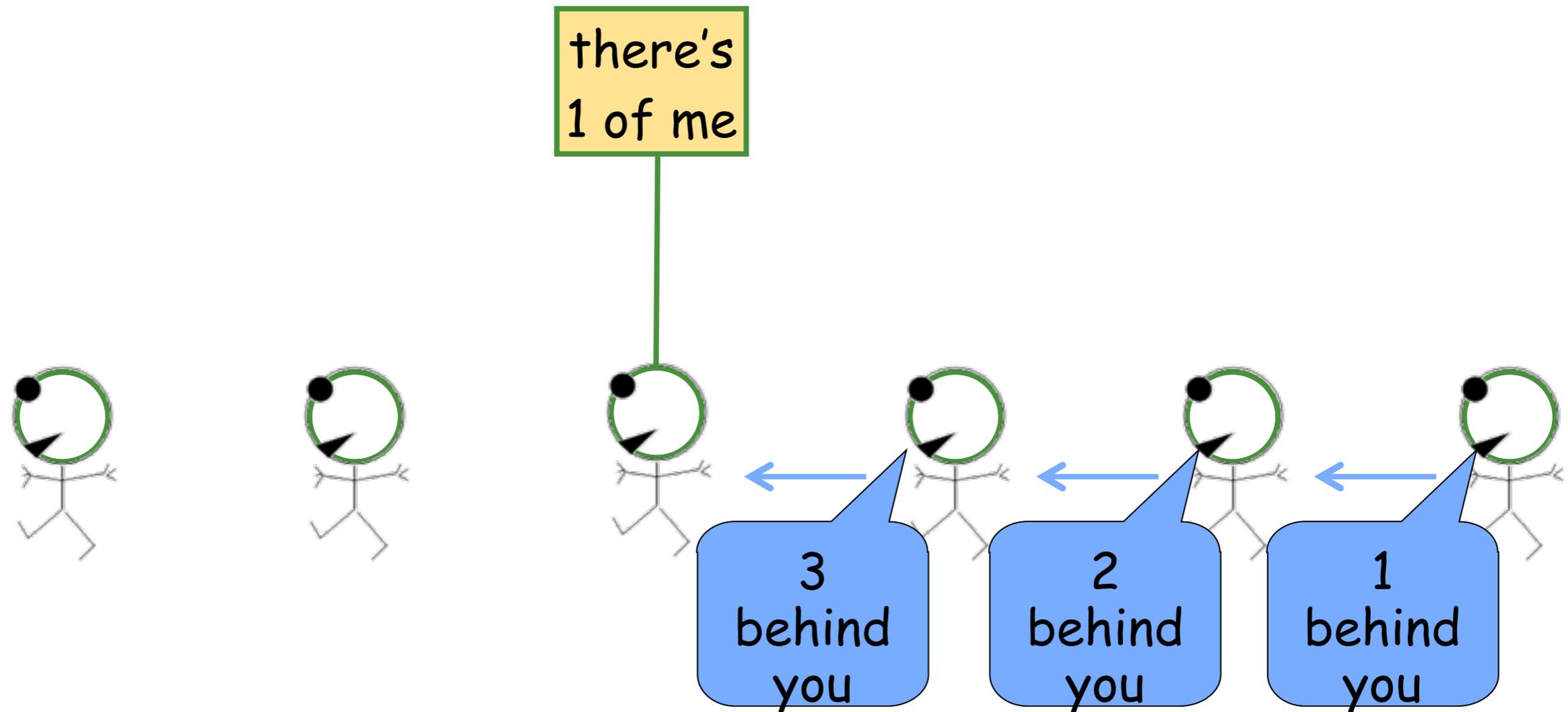
*Count the soldiers*



adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

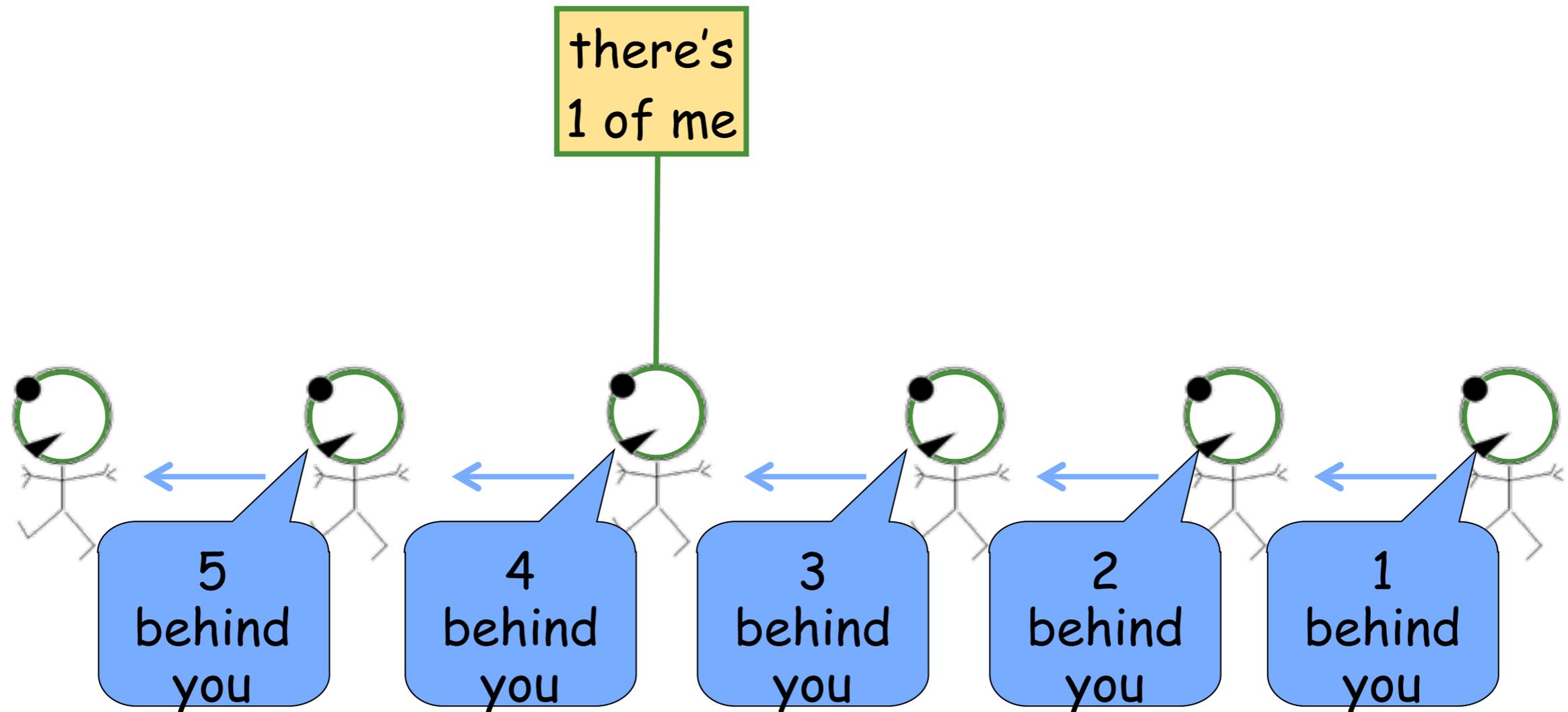
*Count the soldiers*



adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

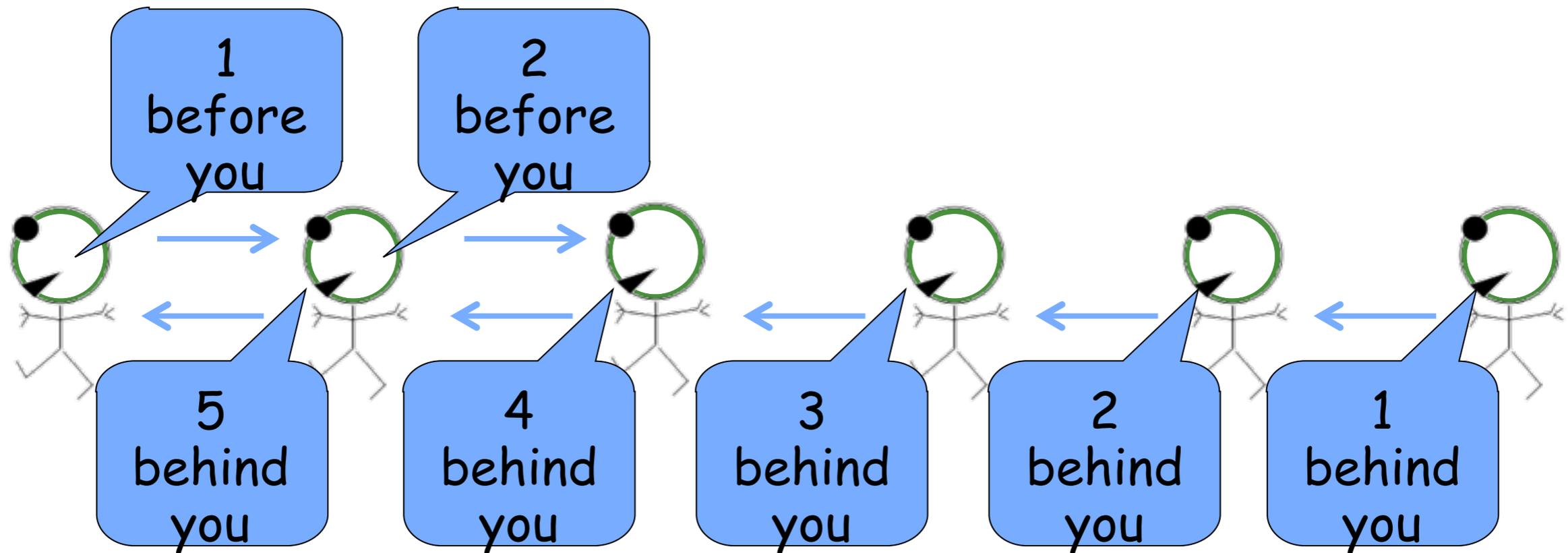
*Count the soldiers*



adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

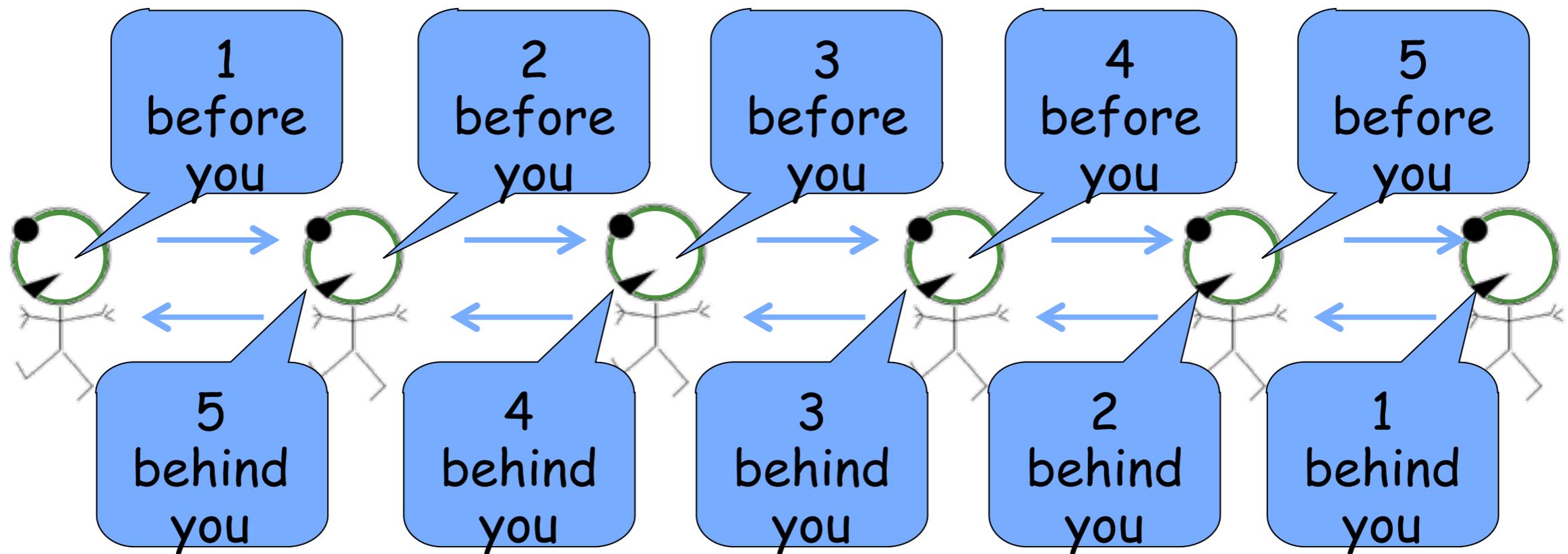
*Count the soldiers*



adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

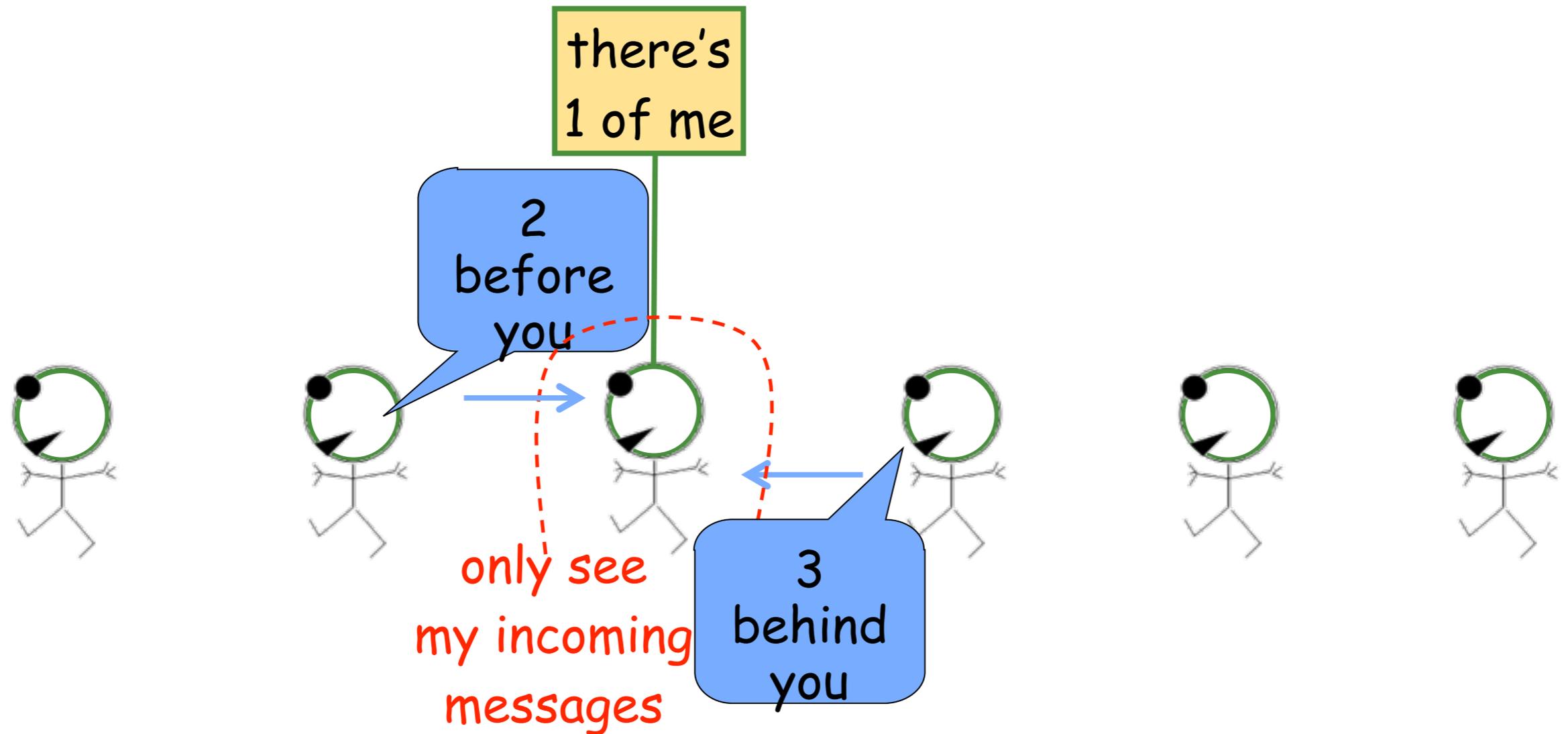
*Count the soldiers*



adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

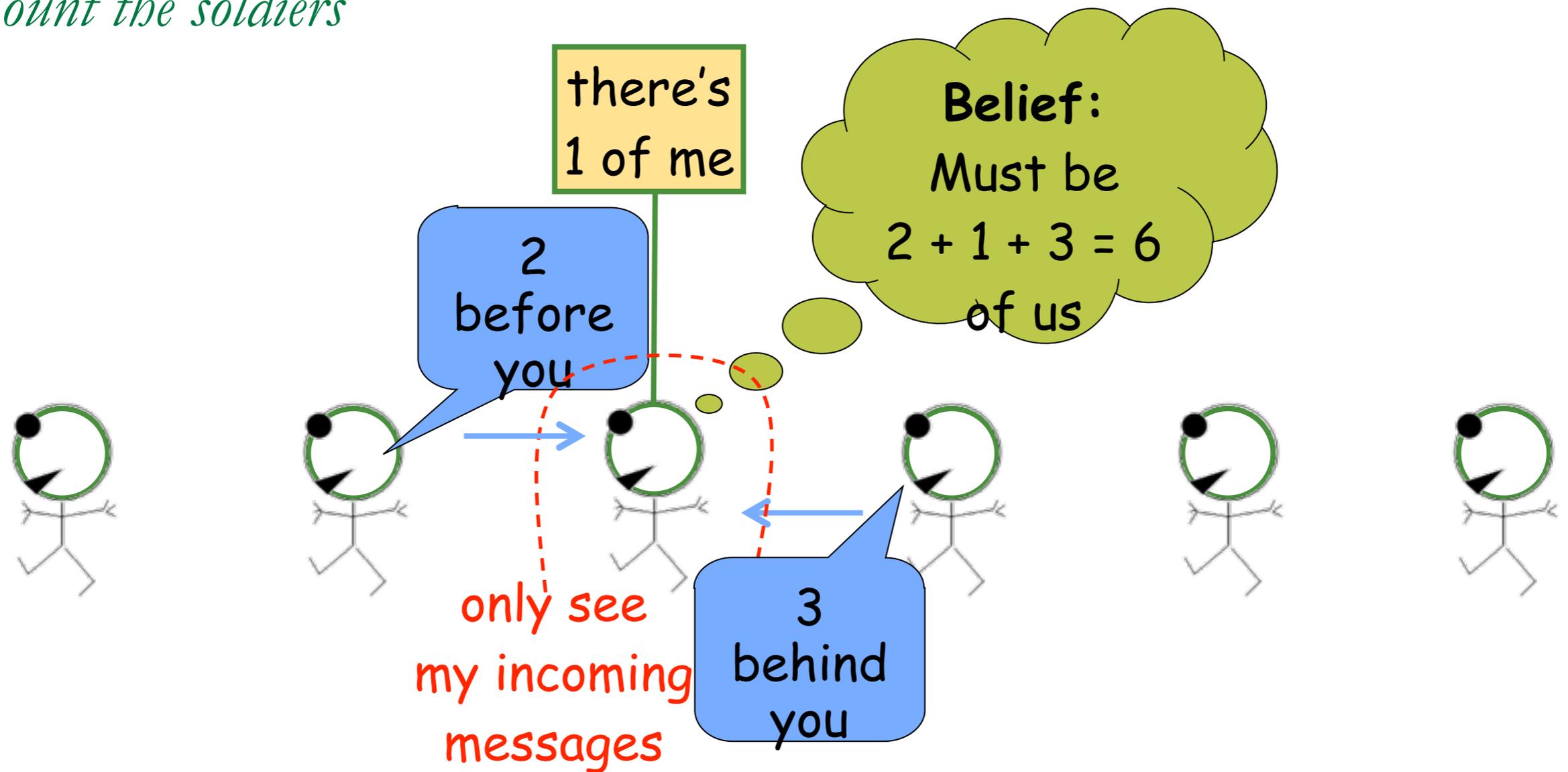
*Count the soldiers*



adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

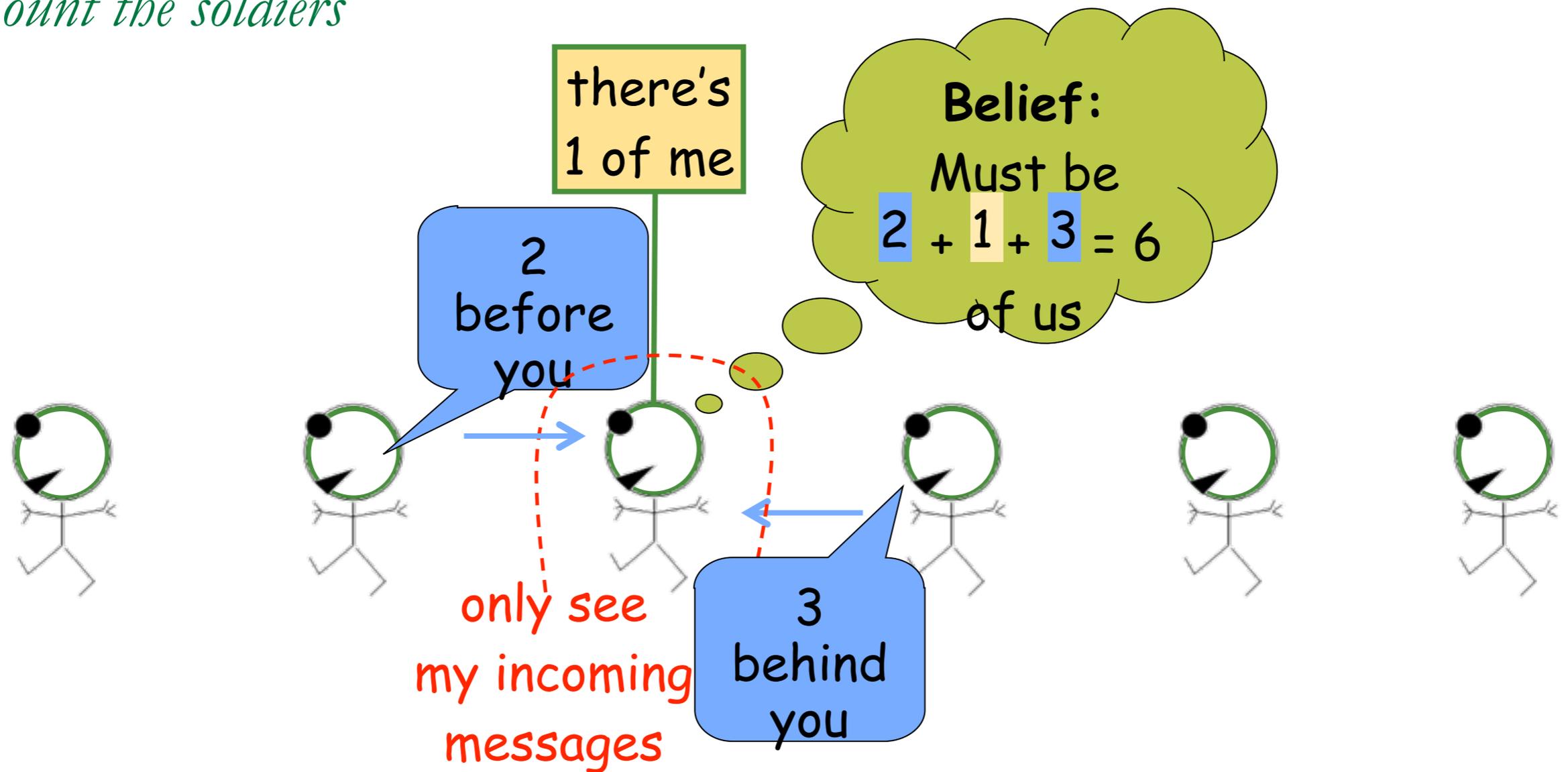
*Count the soldiers*



adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

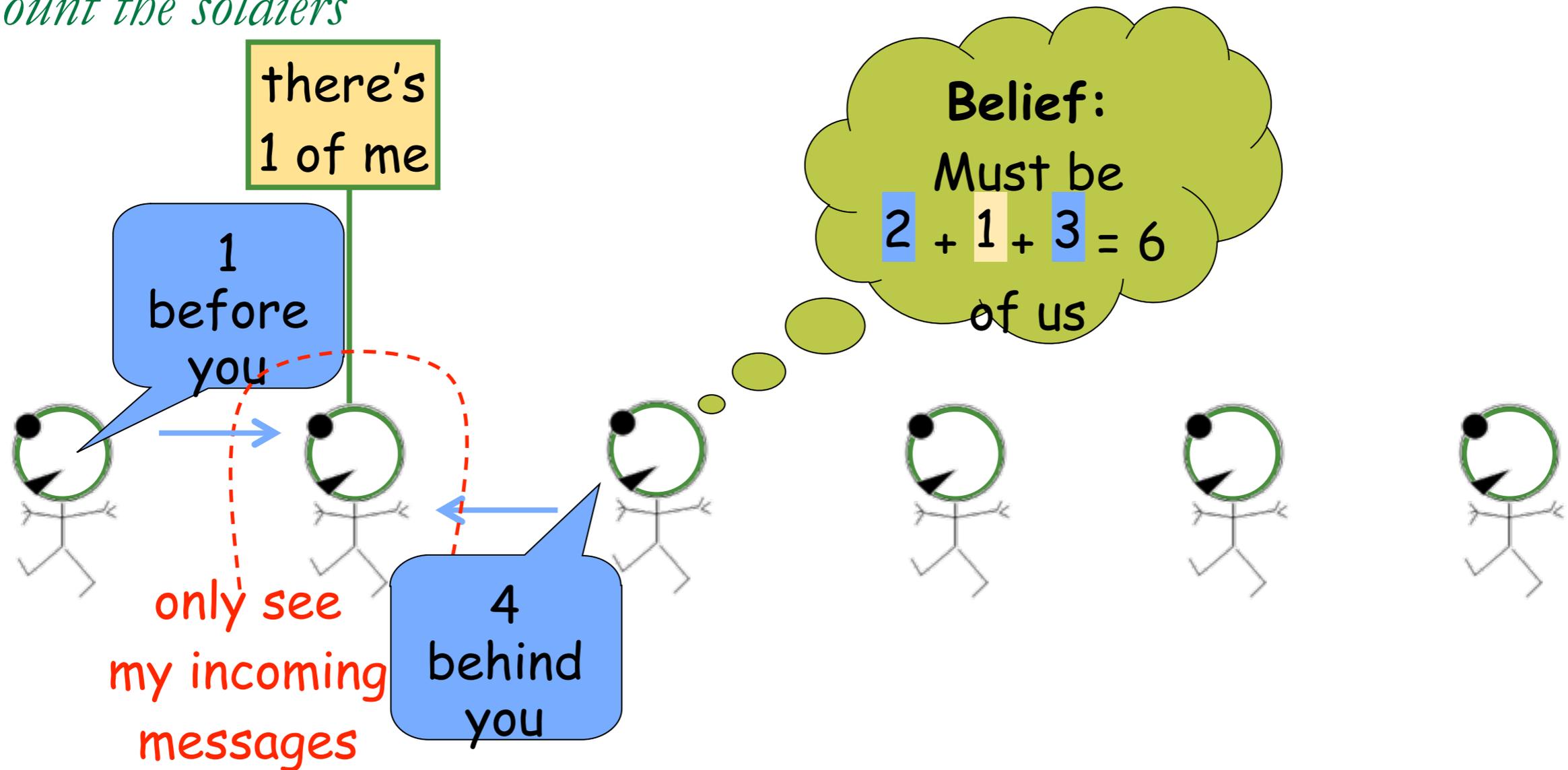
*Count the soldiers*



adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

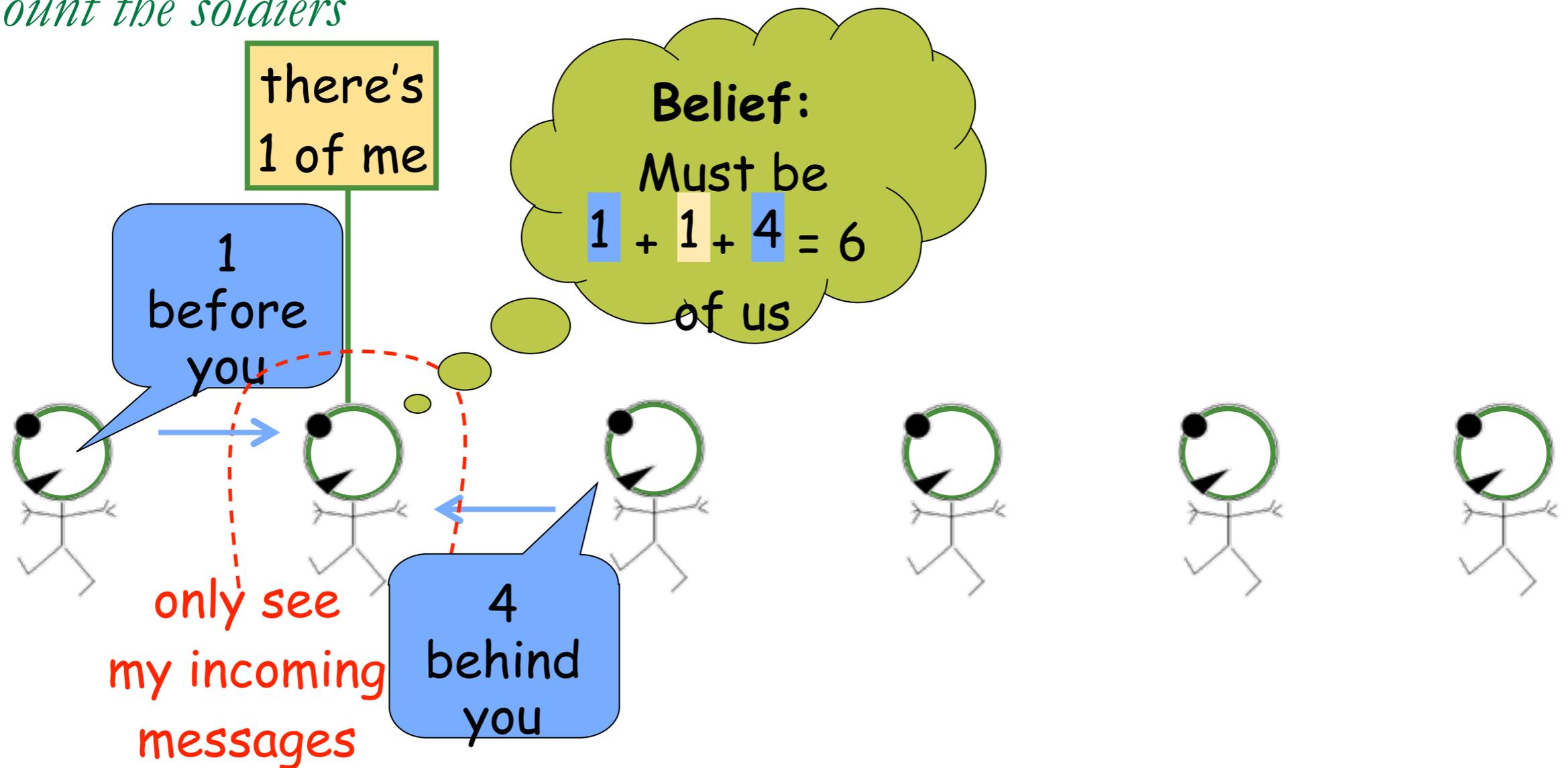
*Count the soldiers*



adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

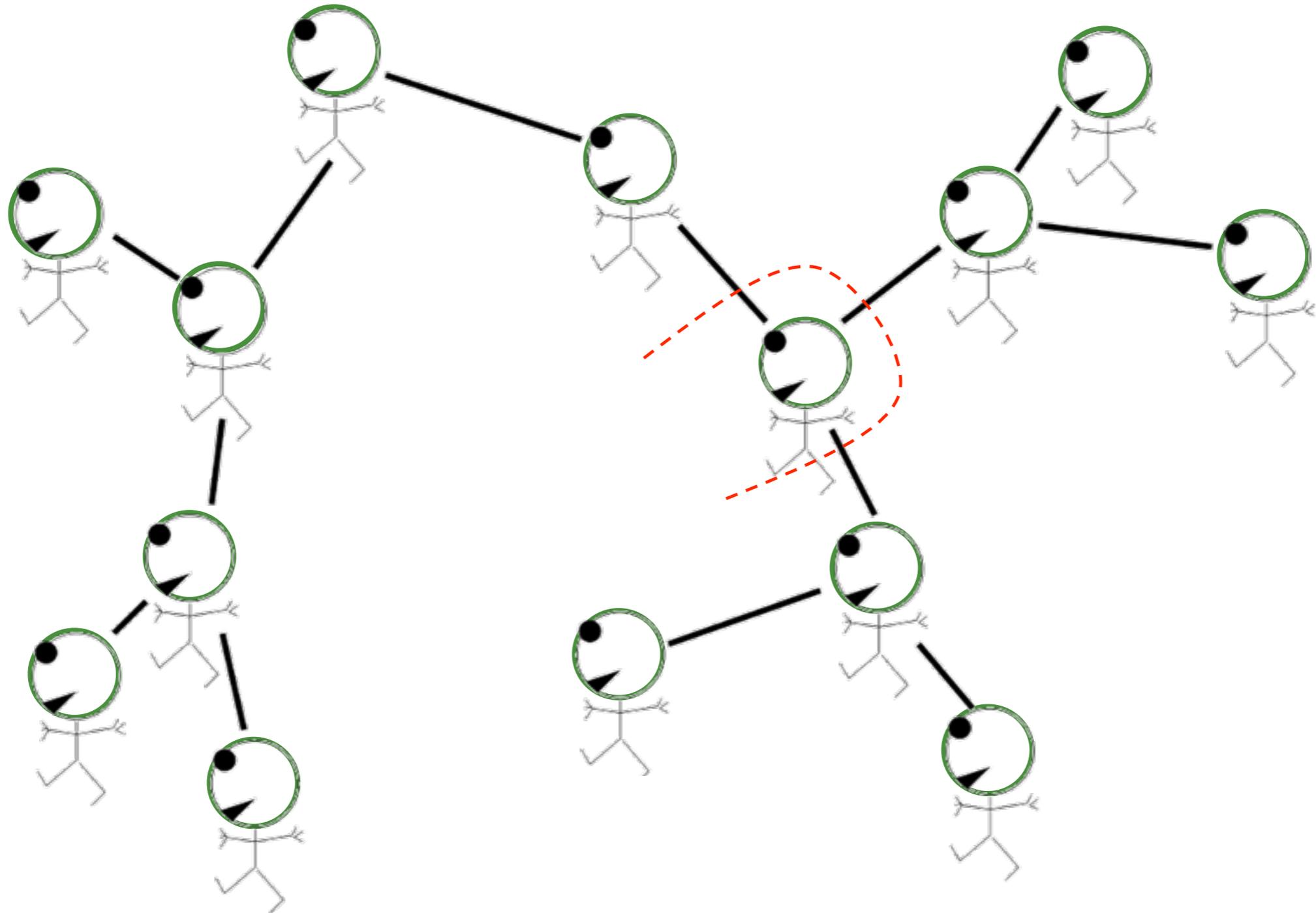
*Count the soldiers*



adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

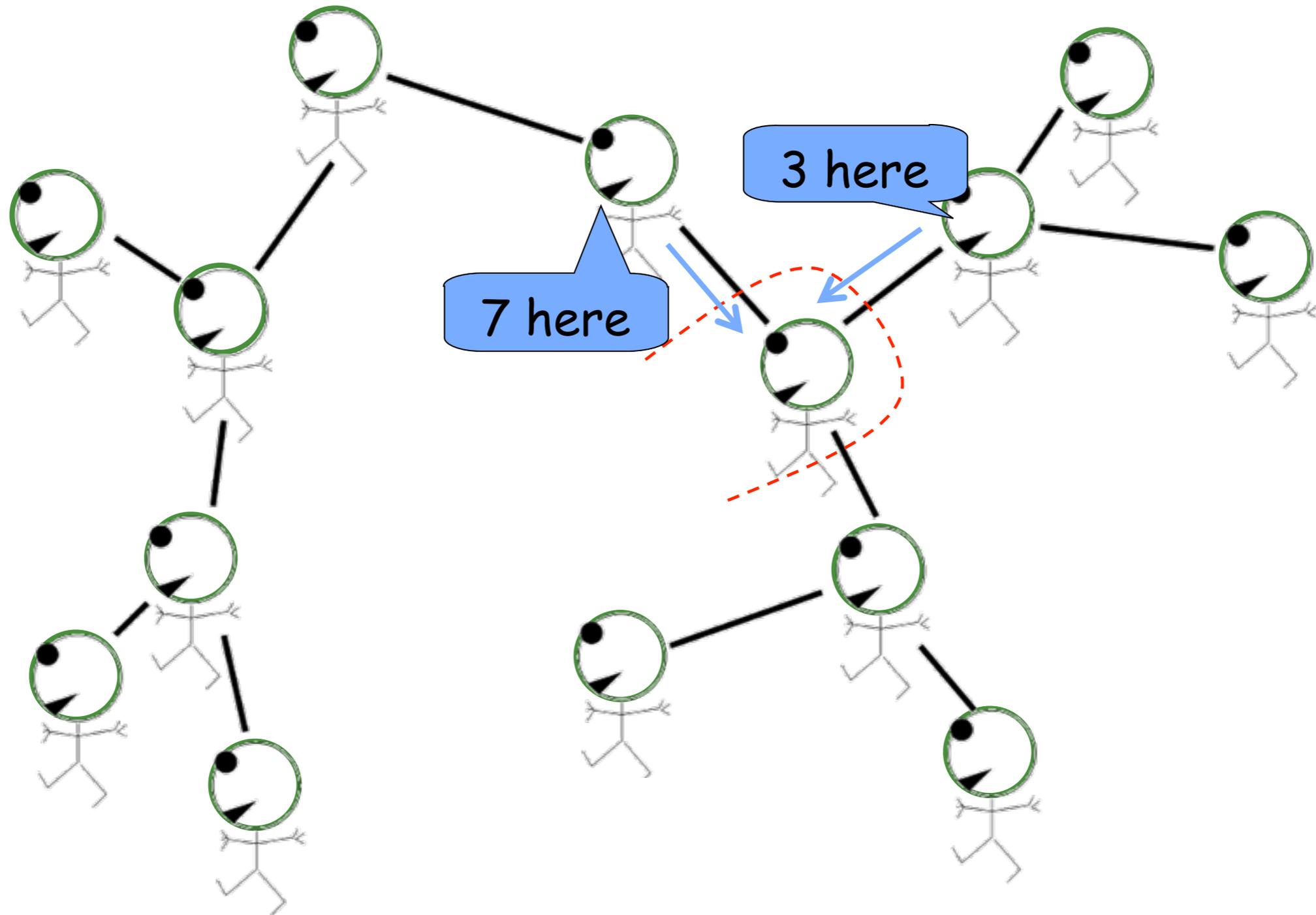
*Each soldier receives reports from all branches of tree*



adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

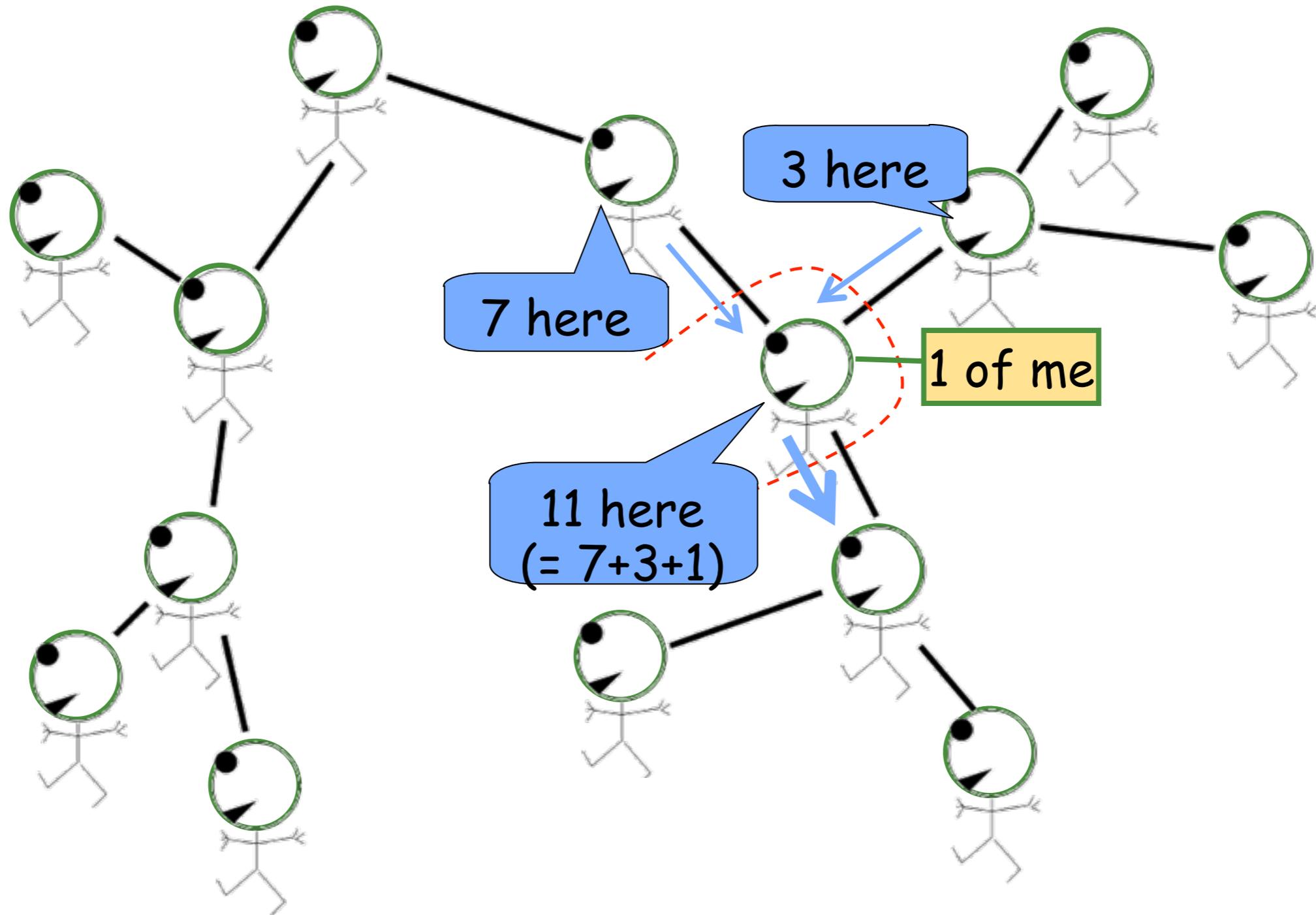
*Each soldier receives reports from all branches of tree*



adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

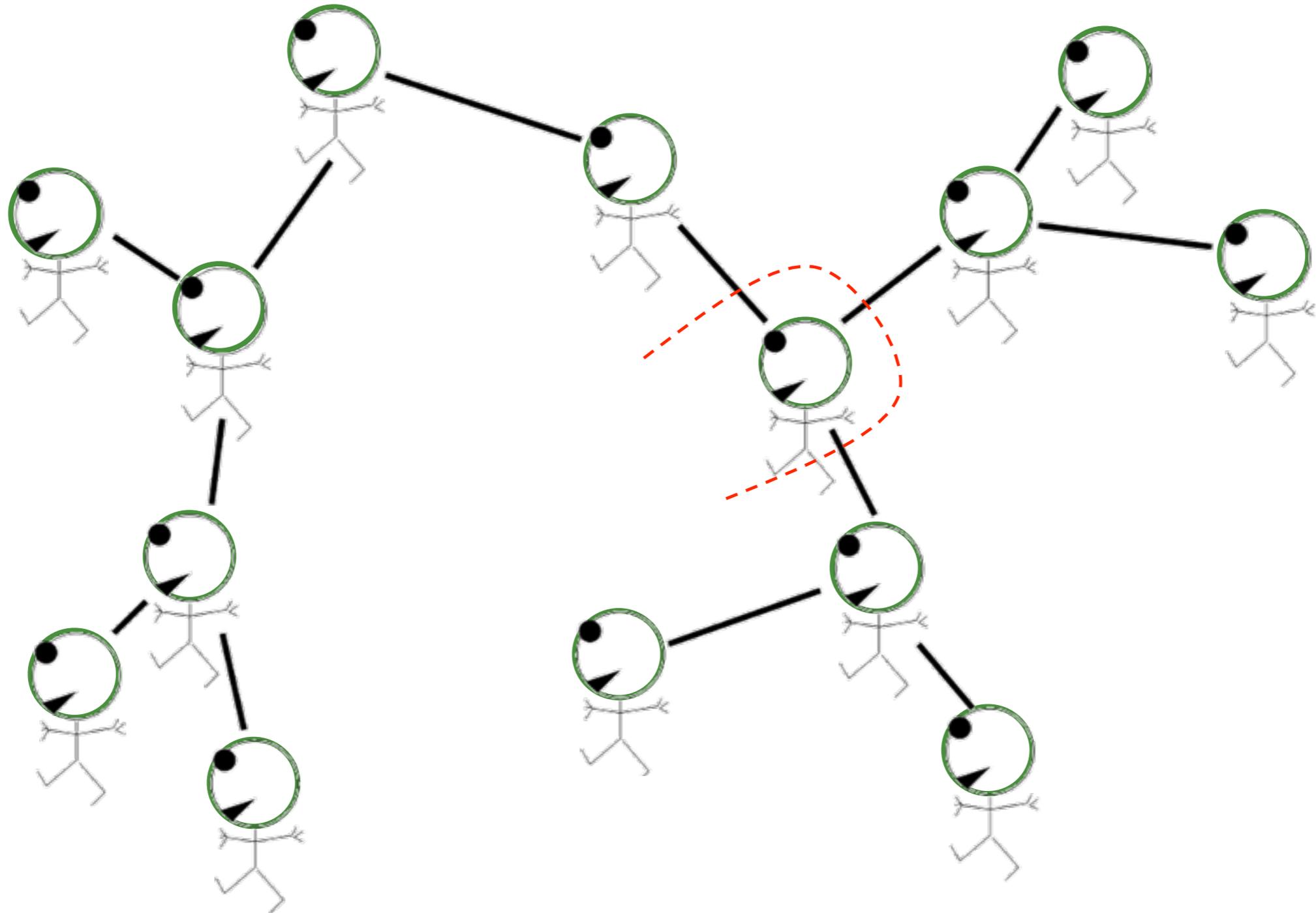
*Each soldier receives reports from all branches of tree*



adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

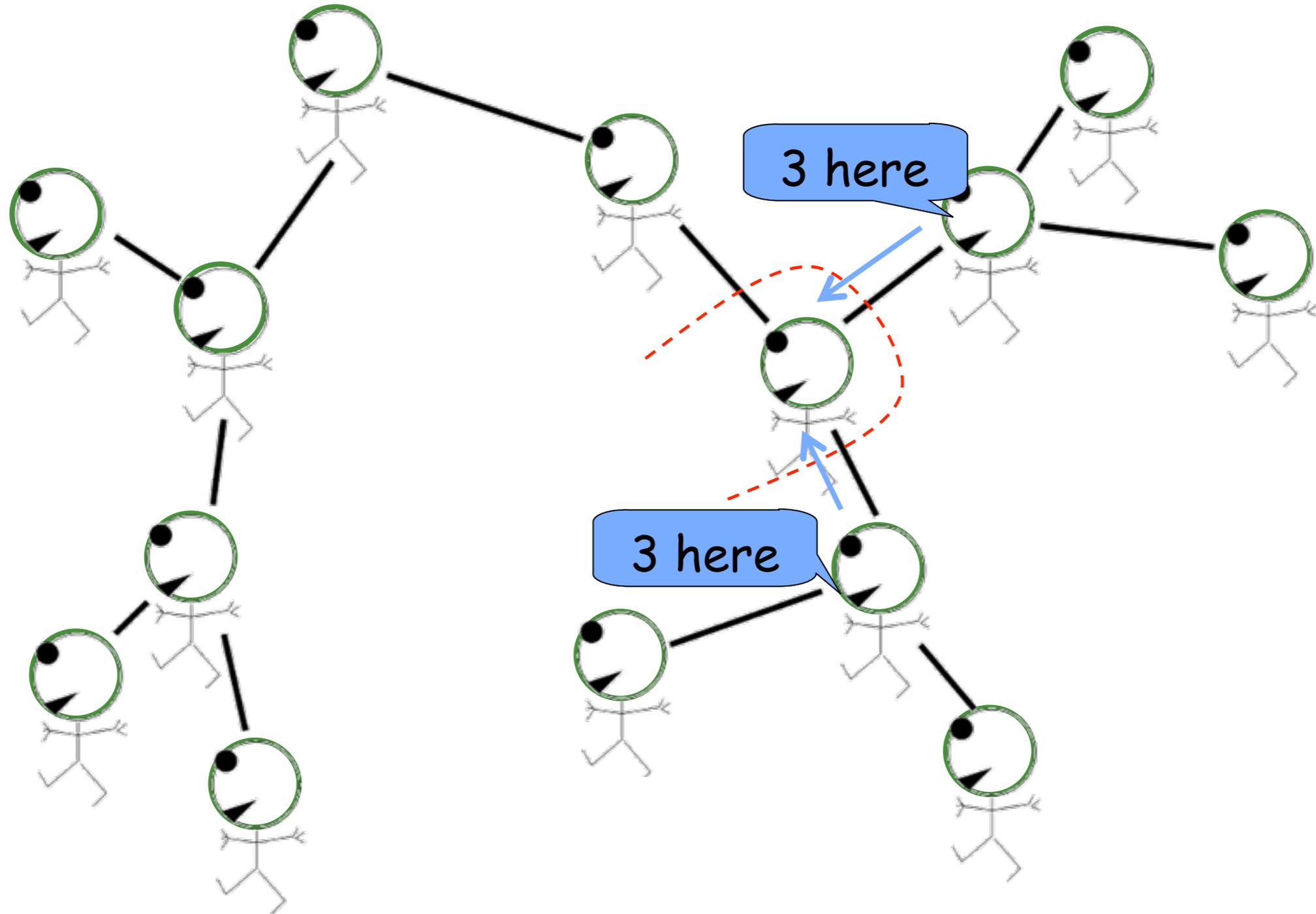
*Each soldier receives reports from all branches of tree*



adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

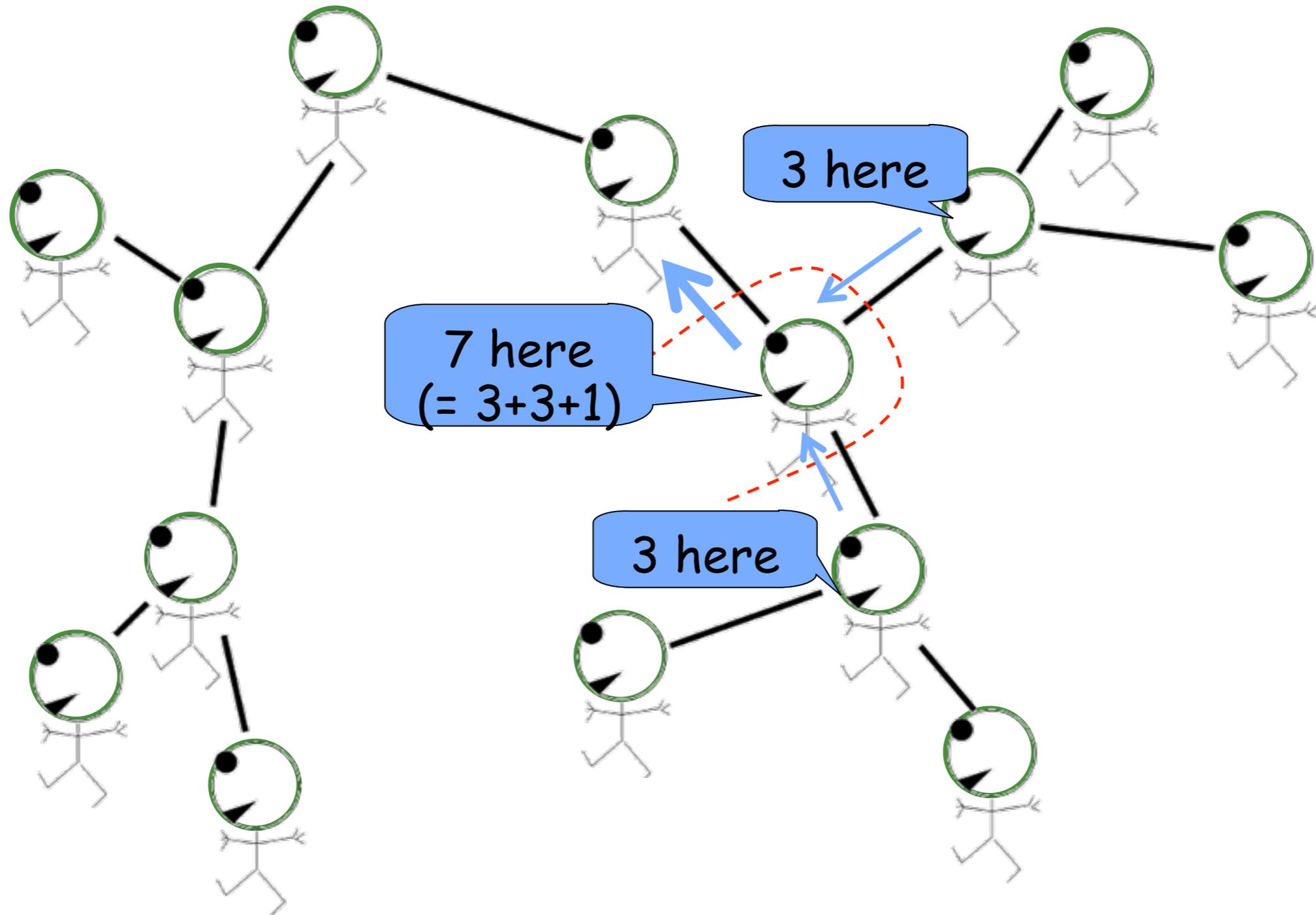
*Each soldier receives reports from all branches of tree*



adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

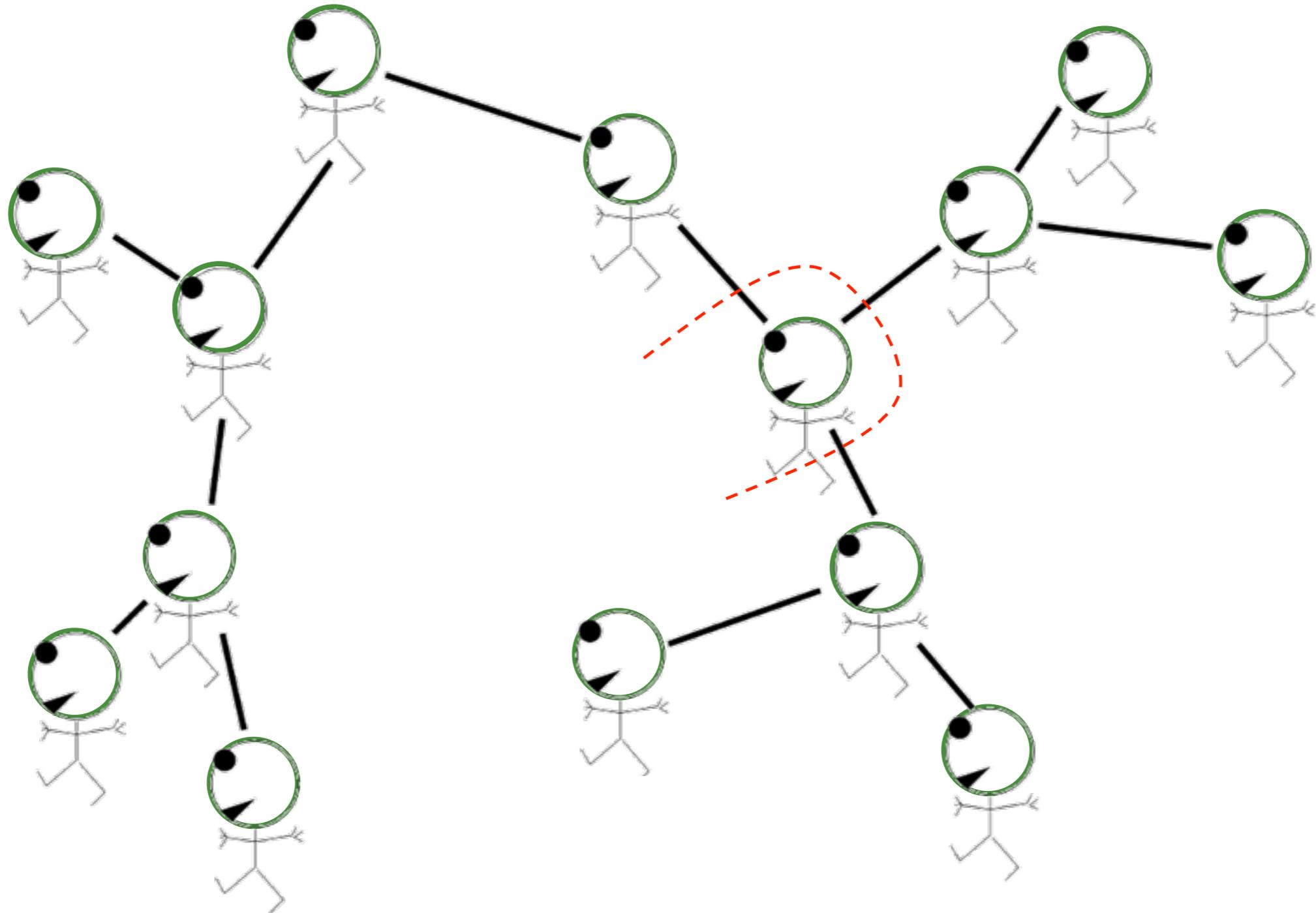
*Each soldier receives reports from all branches of tree*



adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

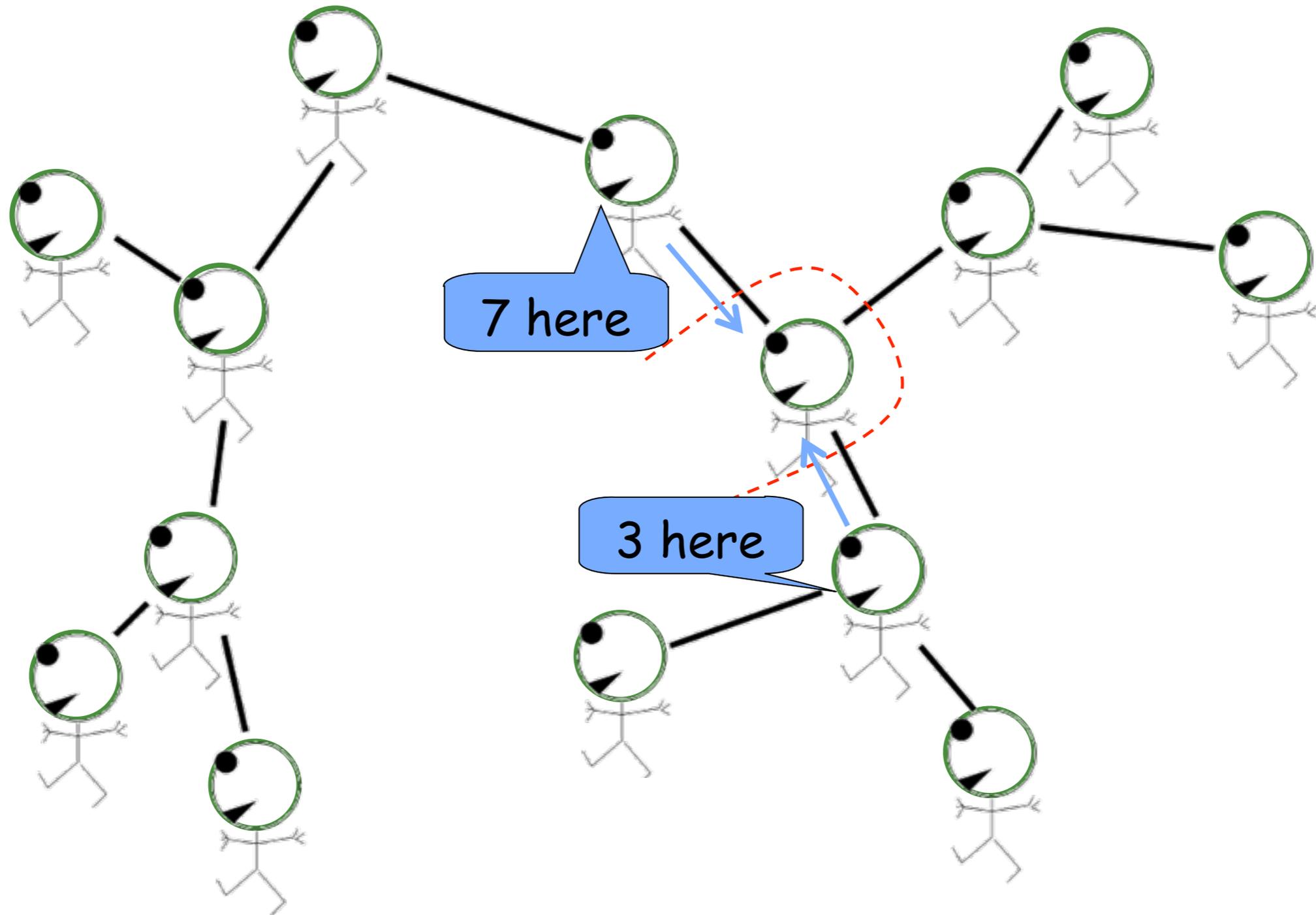
*Each soldier receives reports from all branches of tree*



adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

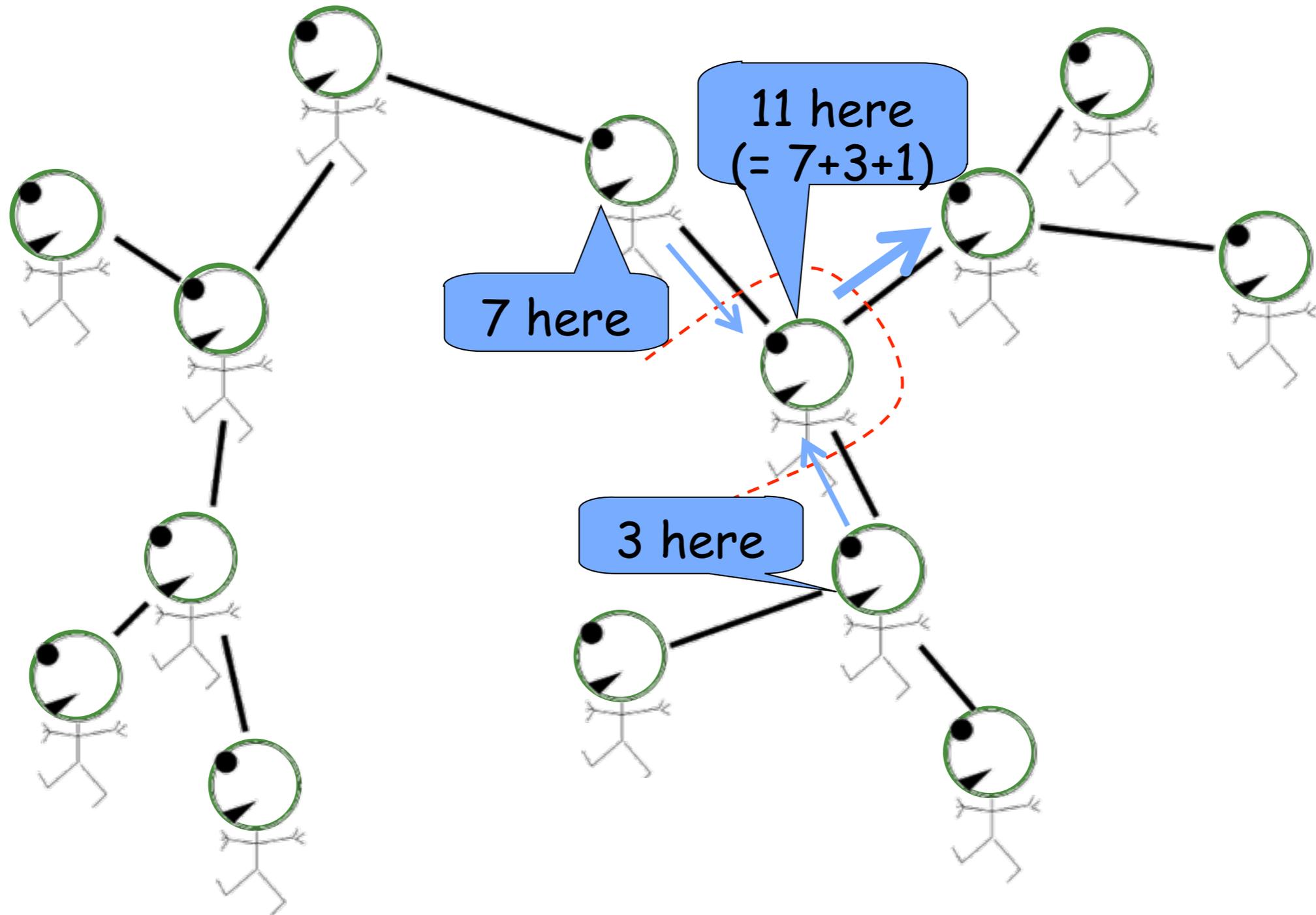
*Each soldier receives reports from all branches of tree*



adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

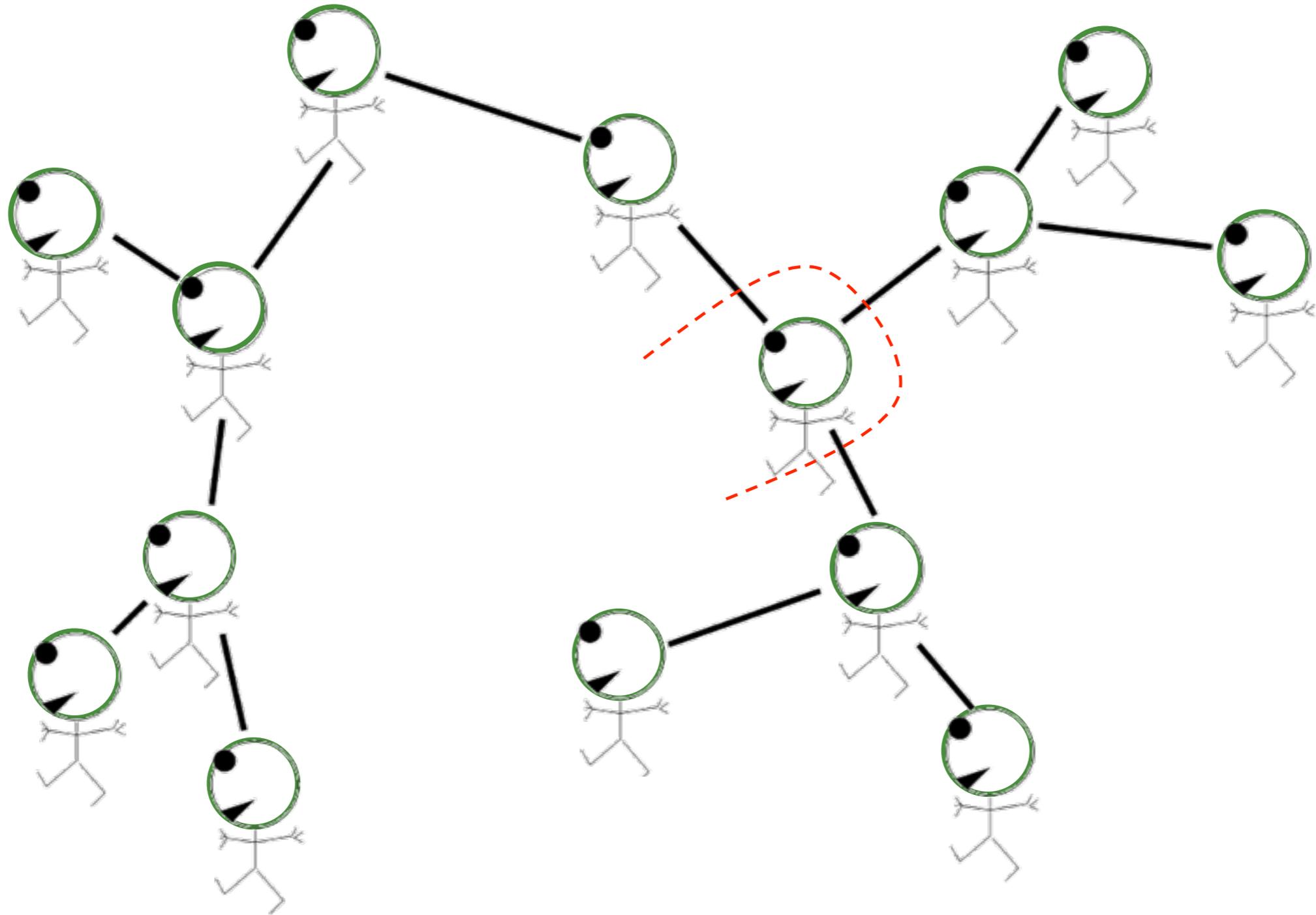
*Each soldier receives reports from all branches of tree*



adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

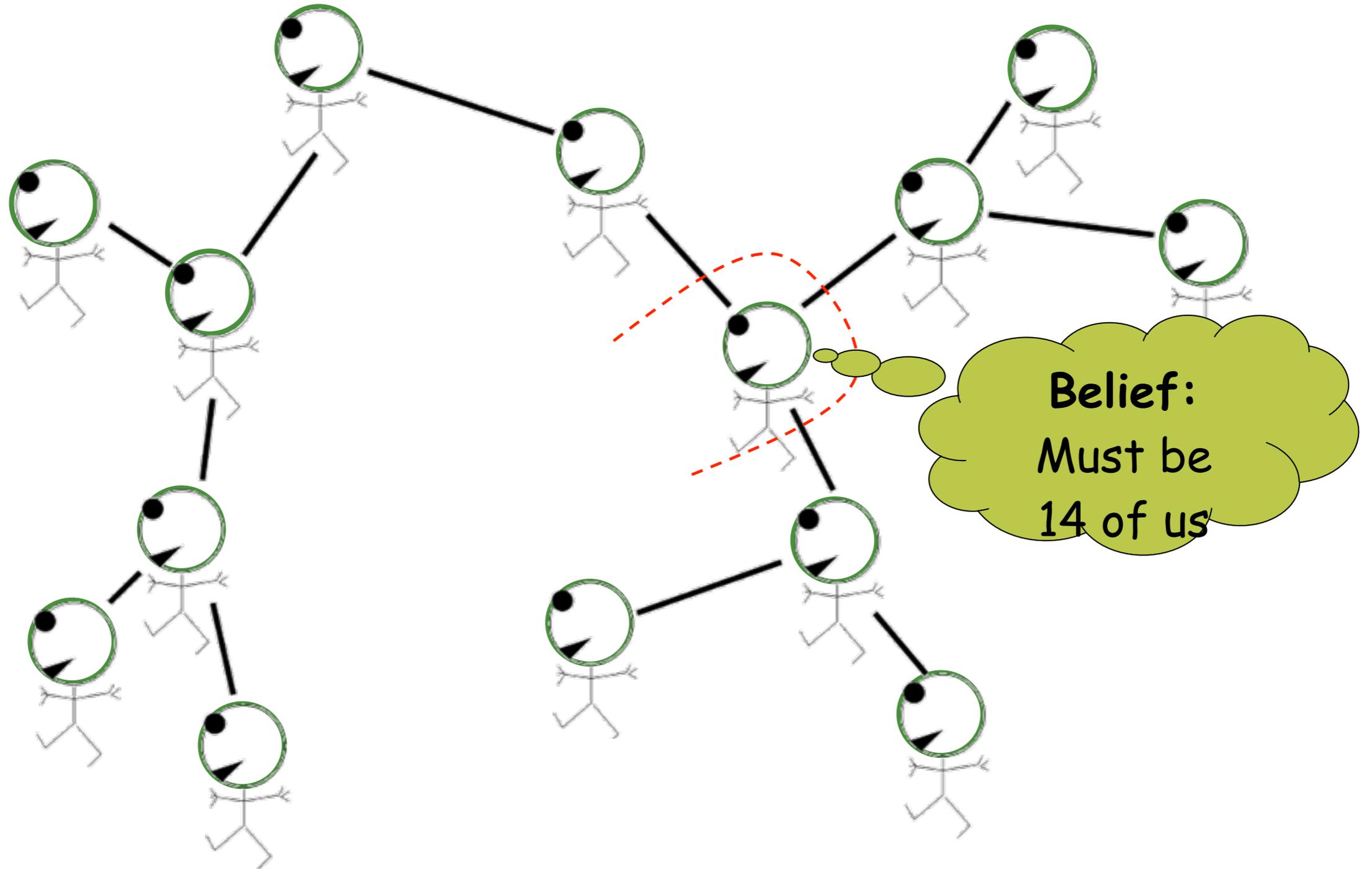
*Each soldier receives reports from all branches of tree*



adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

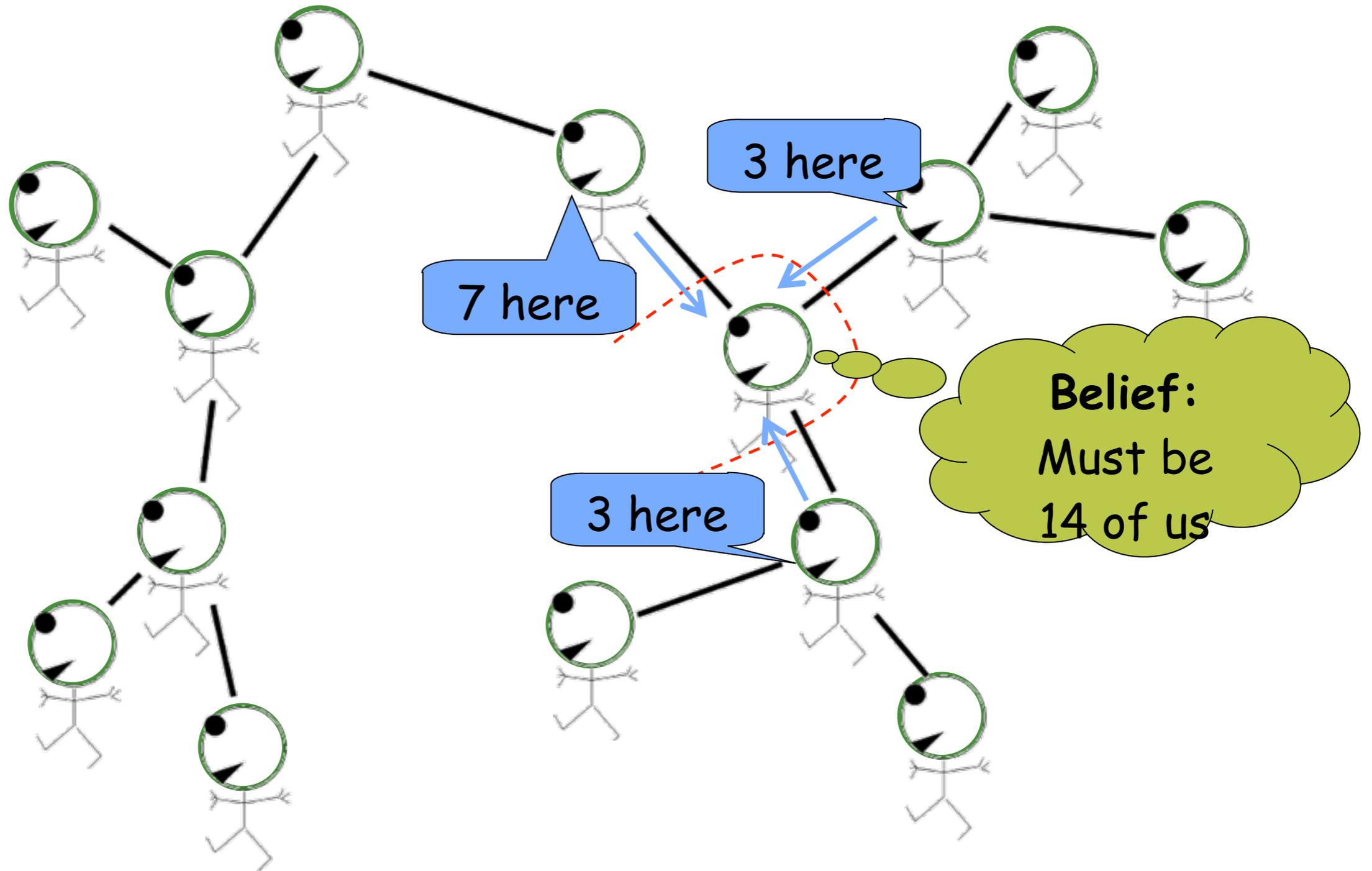
*Each soldier receives reports from all branches of tree*



adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

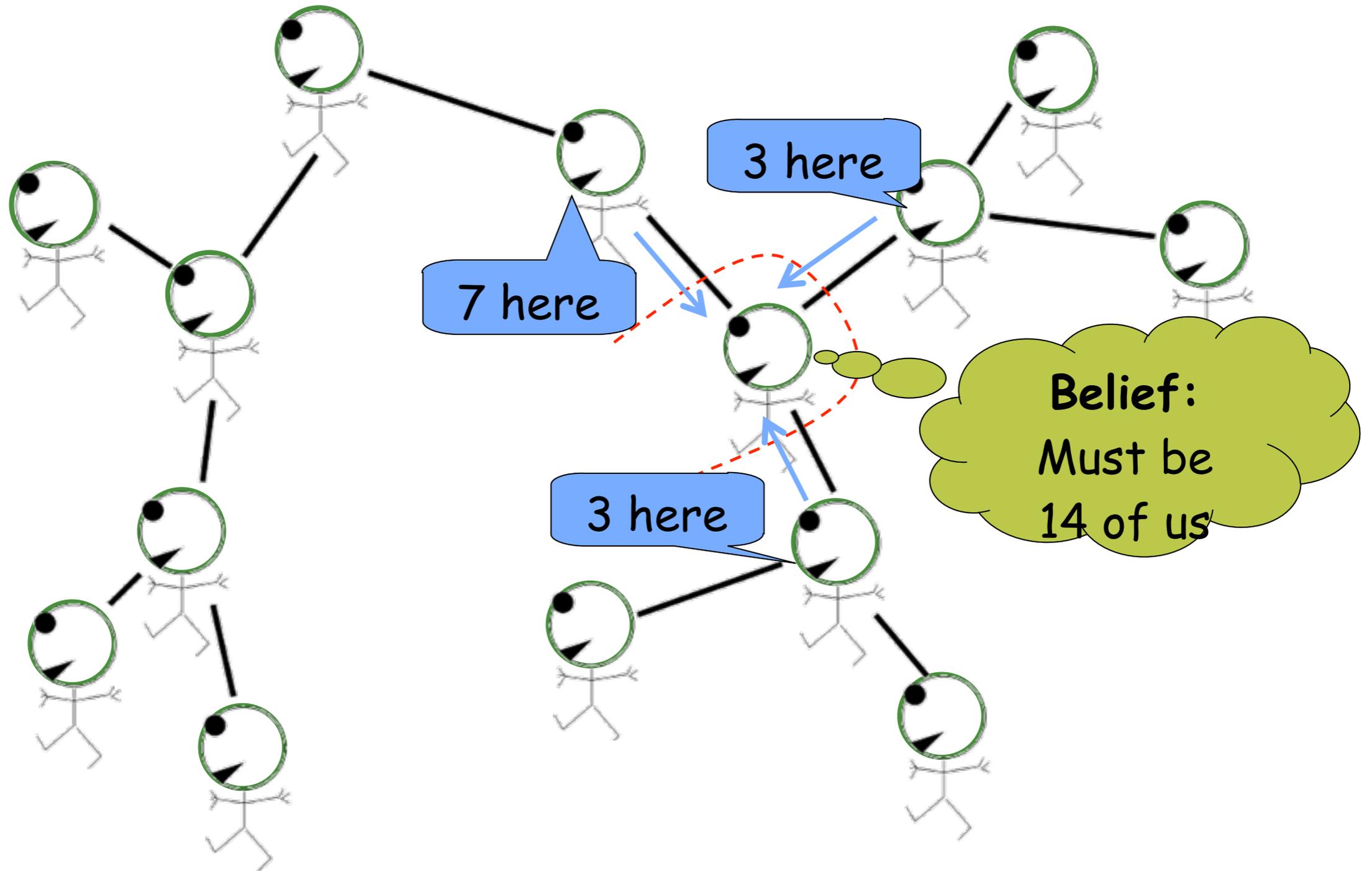
*Each soldier receives reports from all branches of tree*



adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

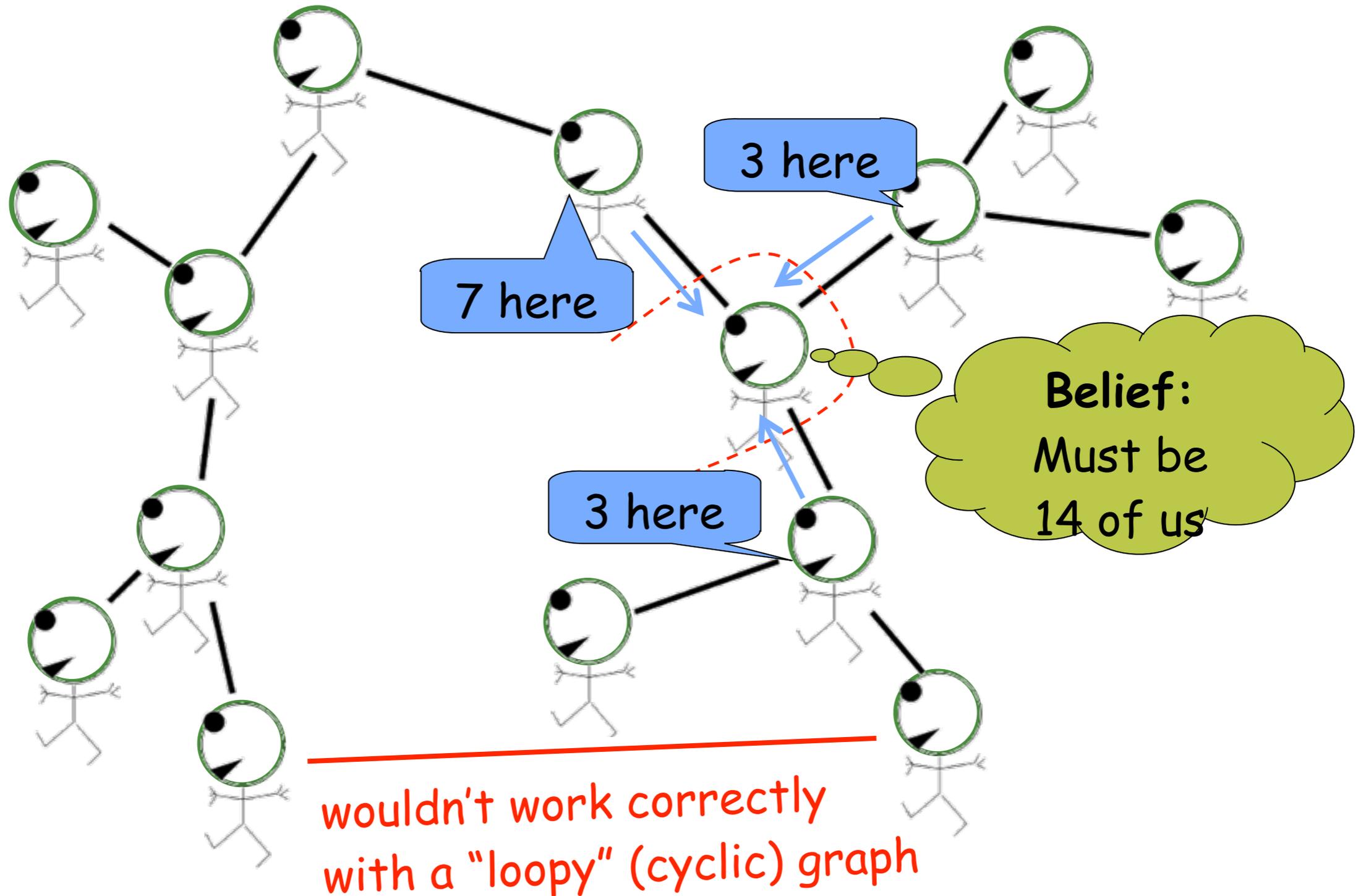
*Each soldier receives reports from all branches of tree*



adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

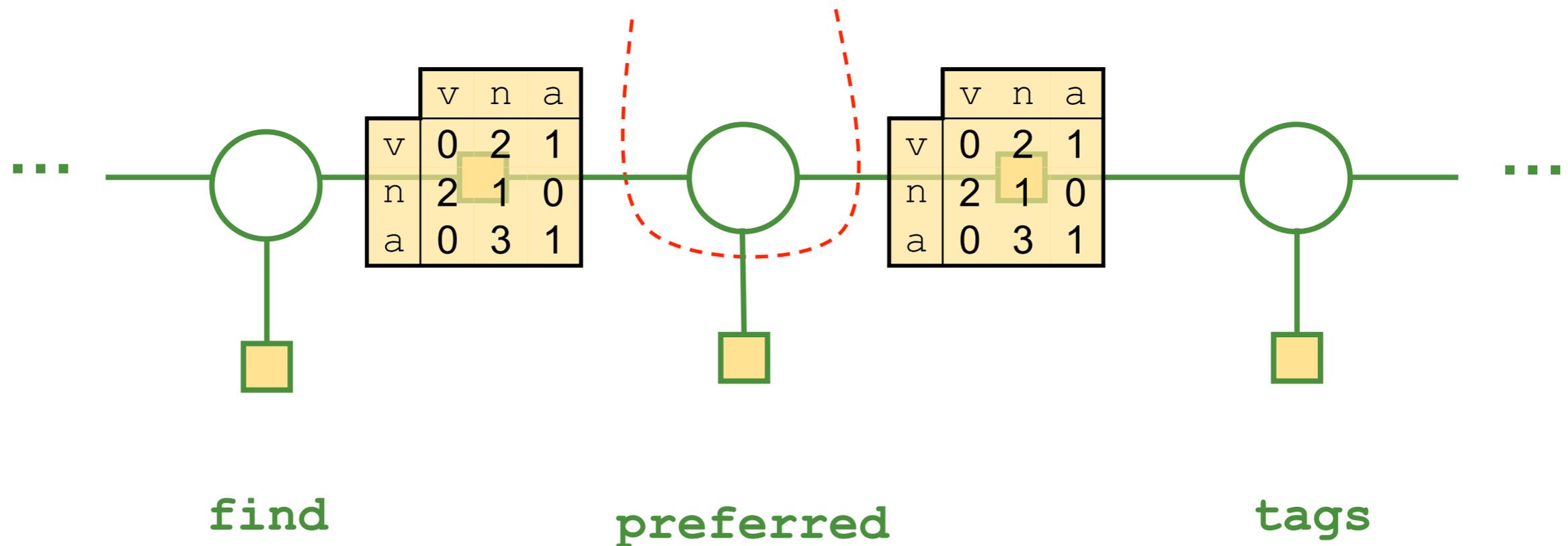
*Each soldier receives reports from all branches of tree*



adapted from MacKay (2003) textbook

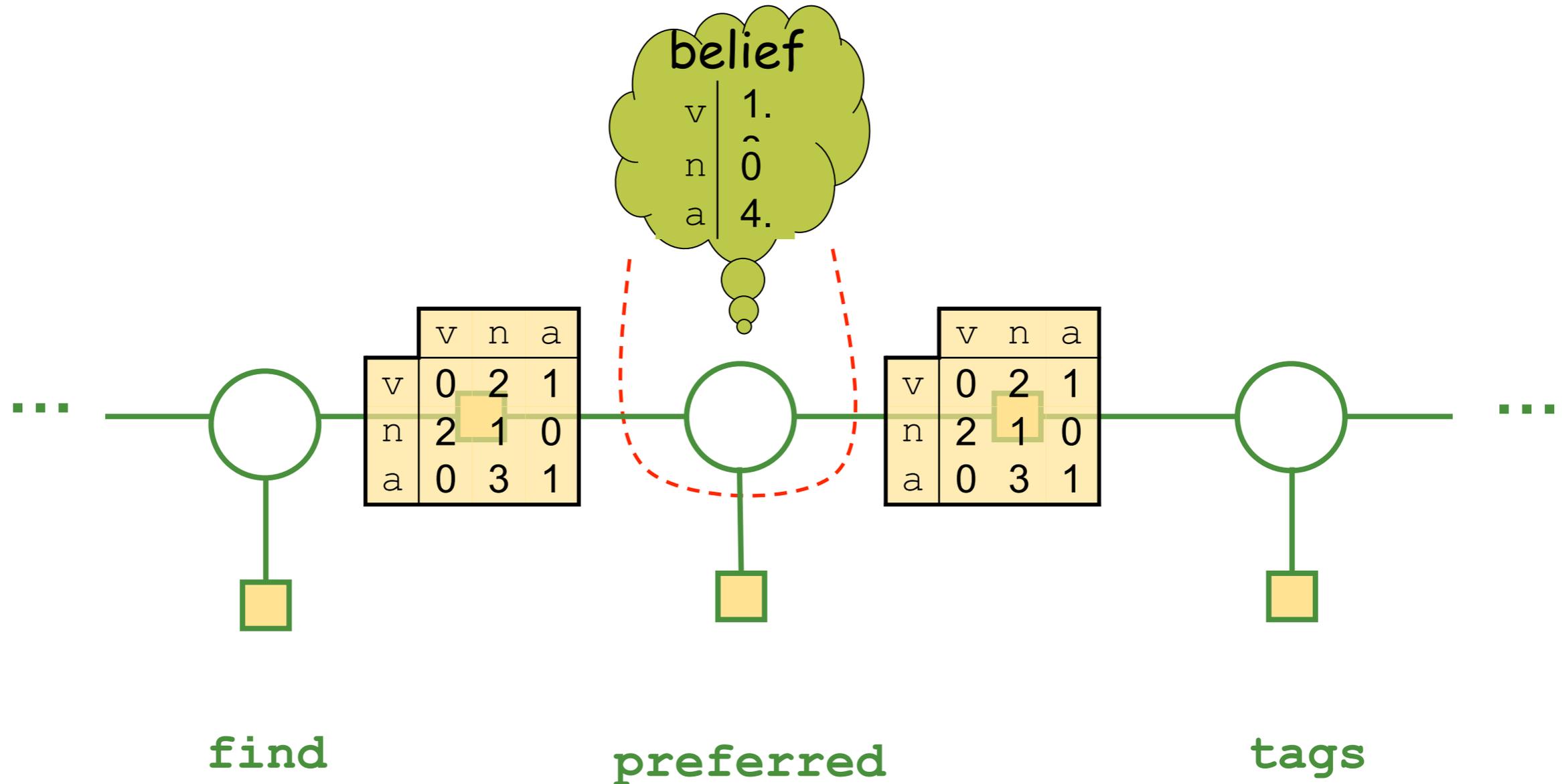
# Great ideas in ML: Forward-Backward

- In the CRF, message passing = forward-backward



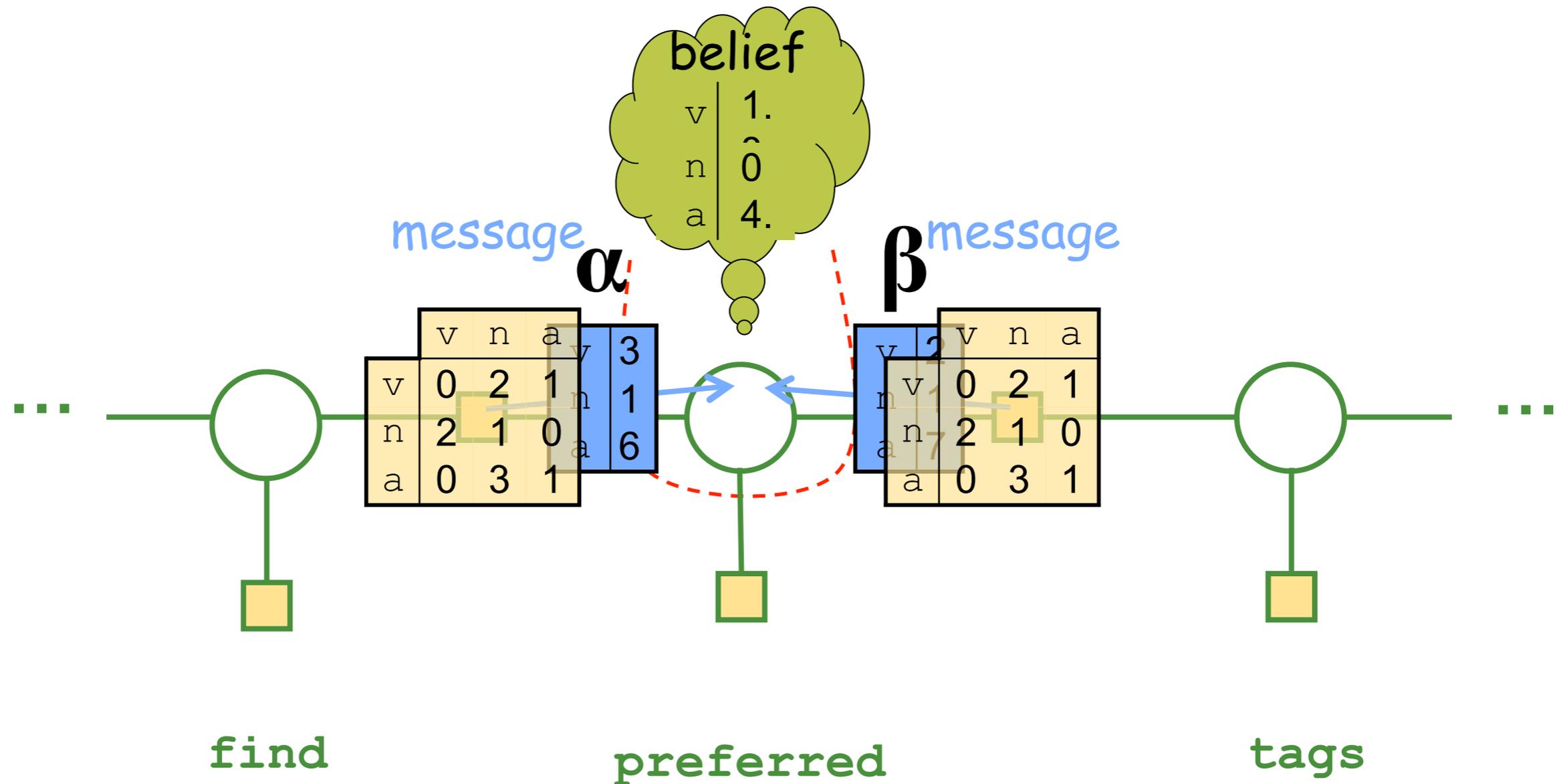
# Great ideas in ML: Forward-Backward

- In the CRF, message passing = forward-backward



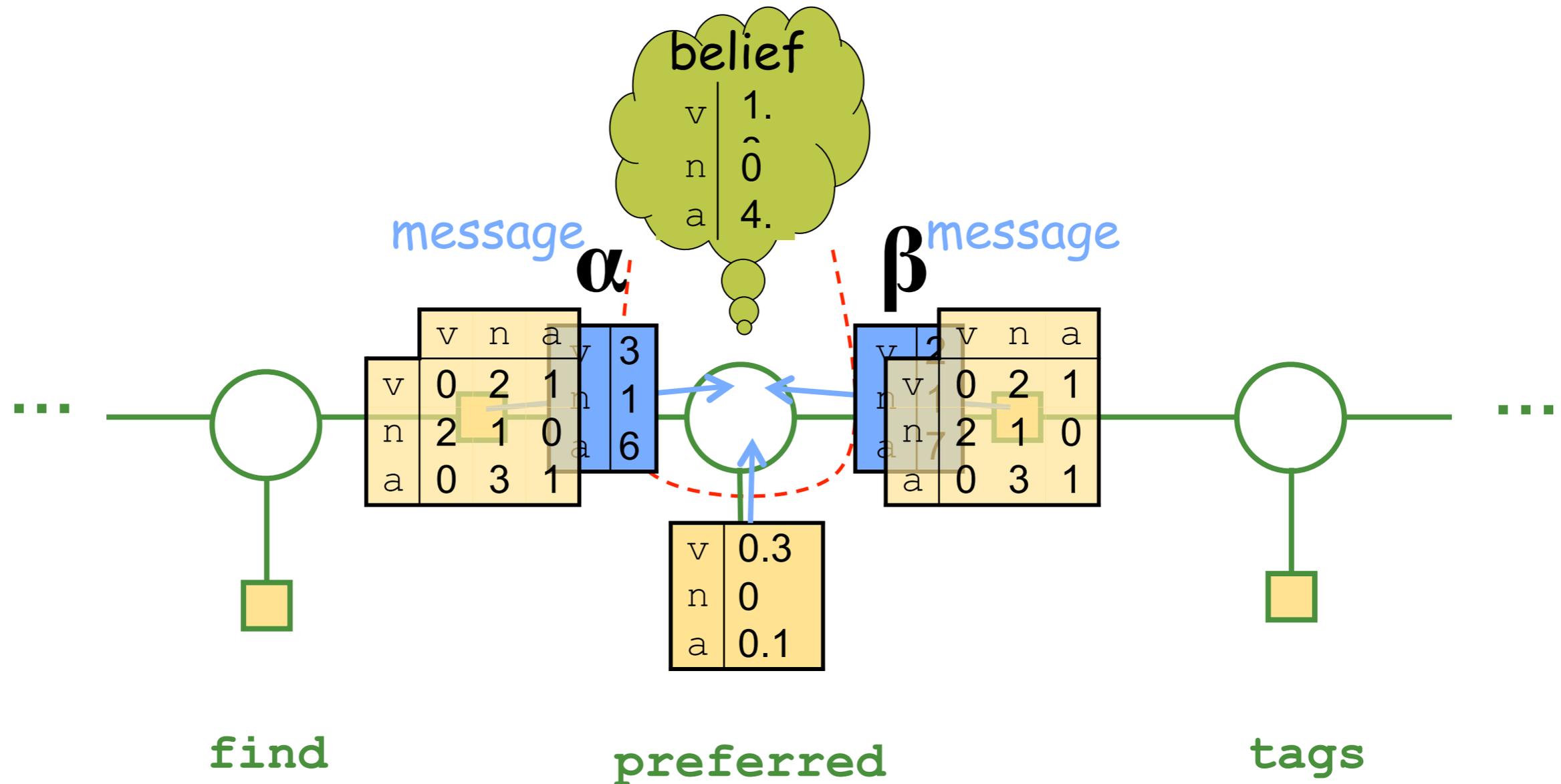
# Great ideas in ML: Forward-Backward

- In the CRF, message passing = forward-backward



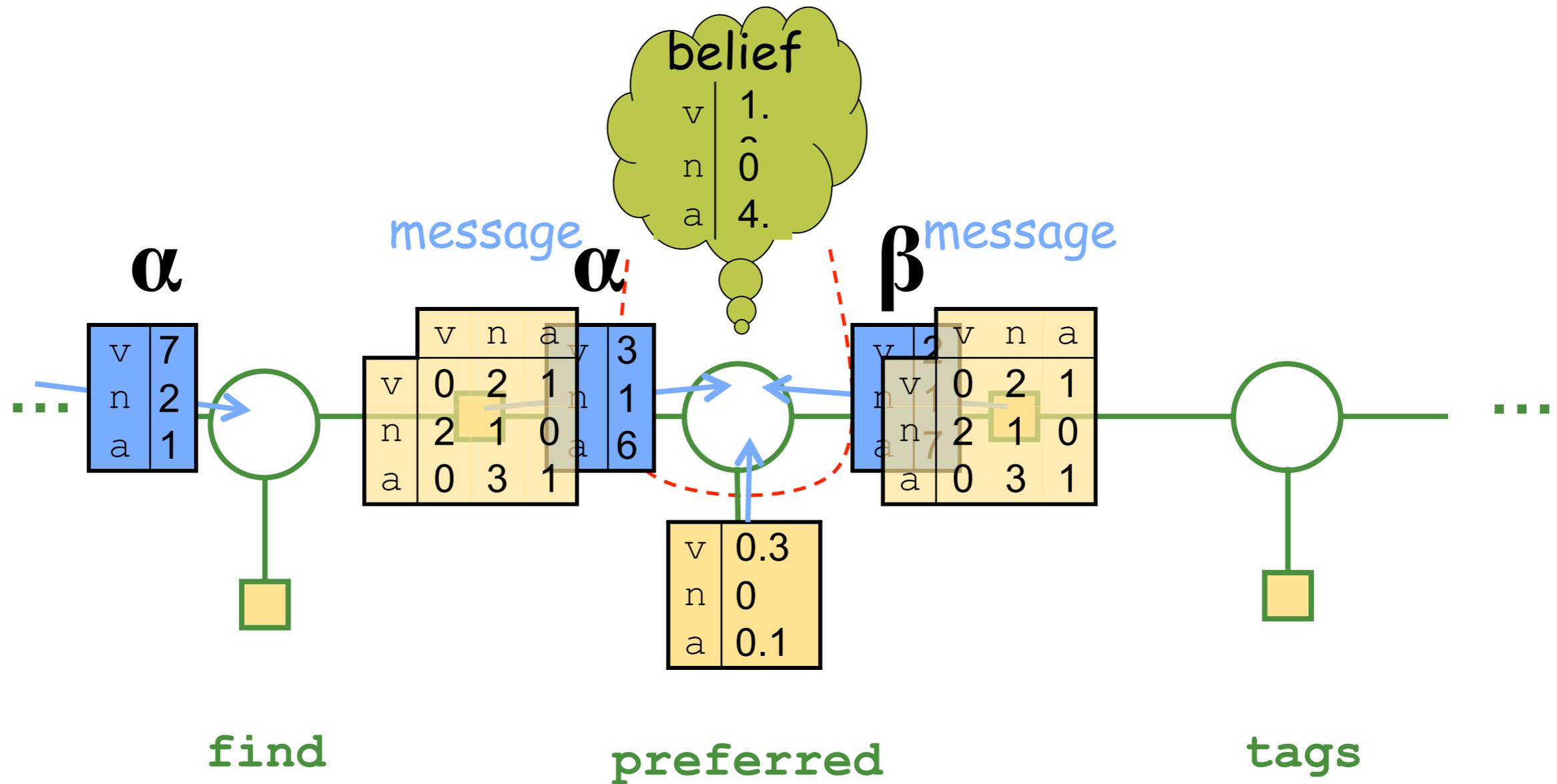
# Great ideas in ML: Forward-Backward

- In the CRF, message passing = forward-backward



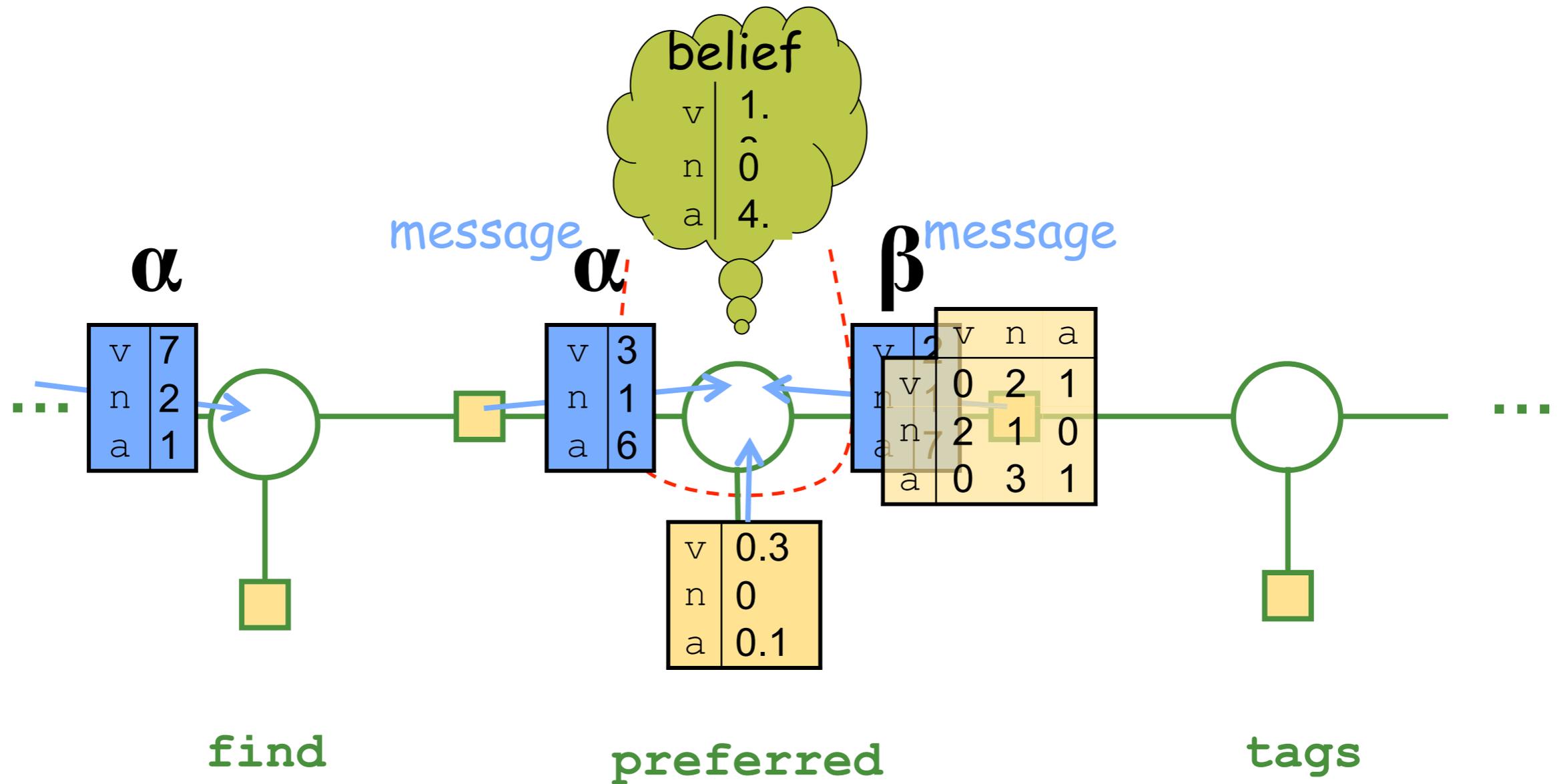
# Great ideas in ML: Forward-Backward

- In the CRF, message passing = forward-backward



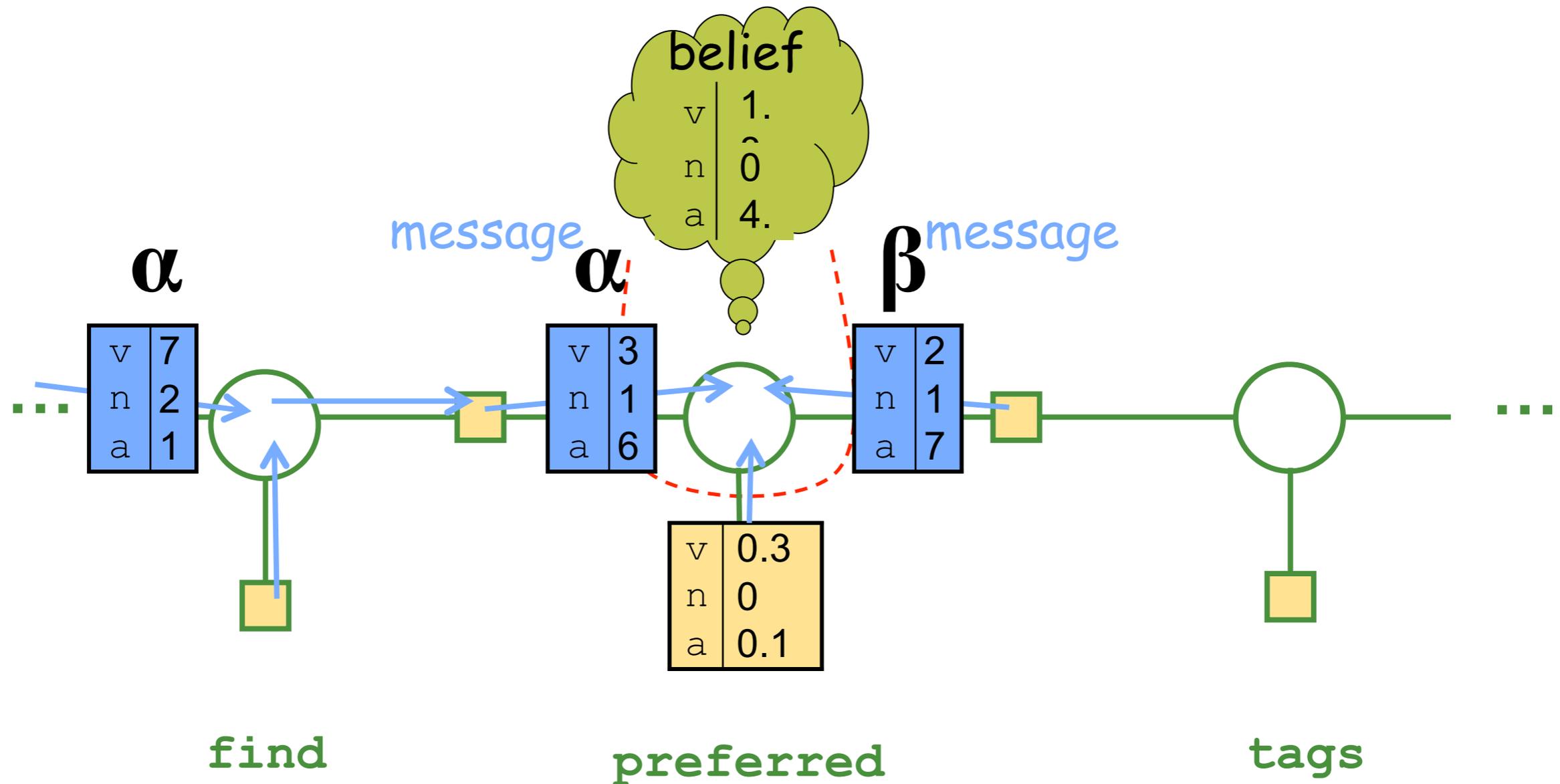
# Great ideas in ML: Forward-Backward

- In the CRF, message passing = forward-backward



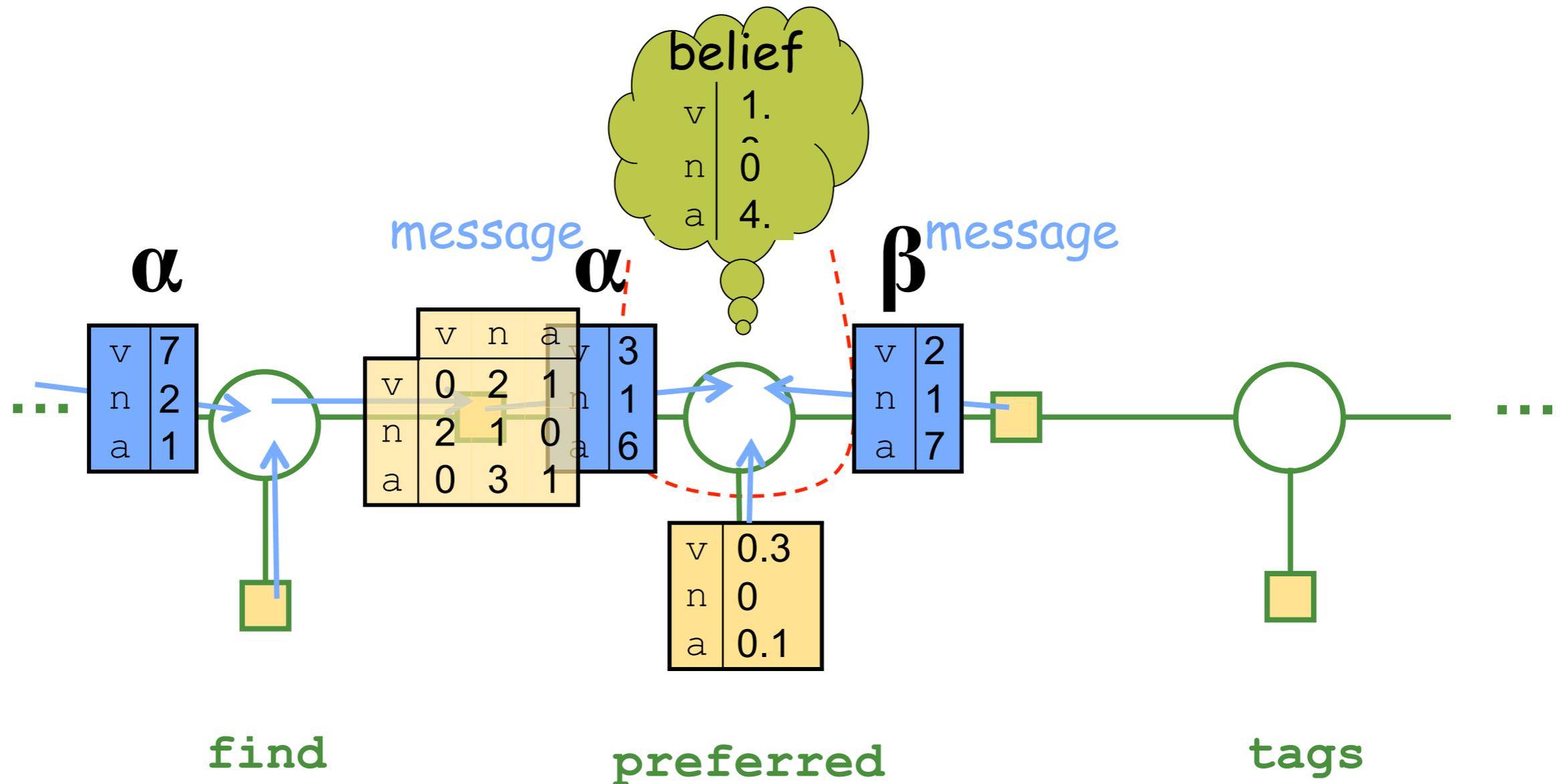
# Great ideas in ML: Forward-Backward

- In the CRF, message passing = forward-backward



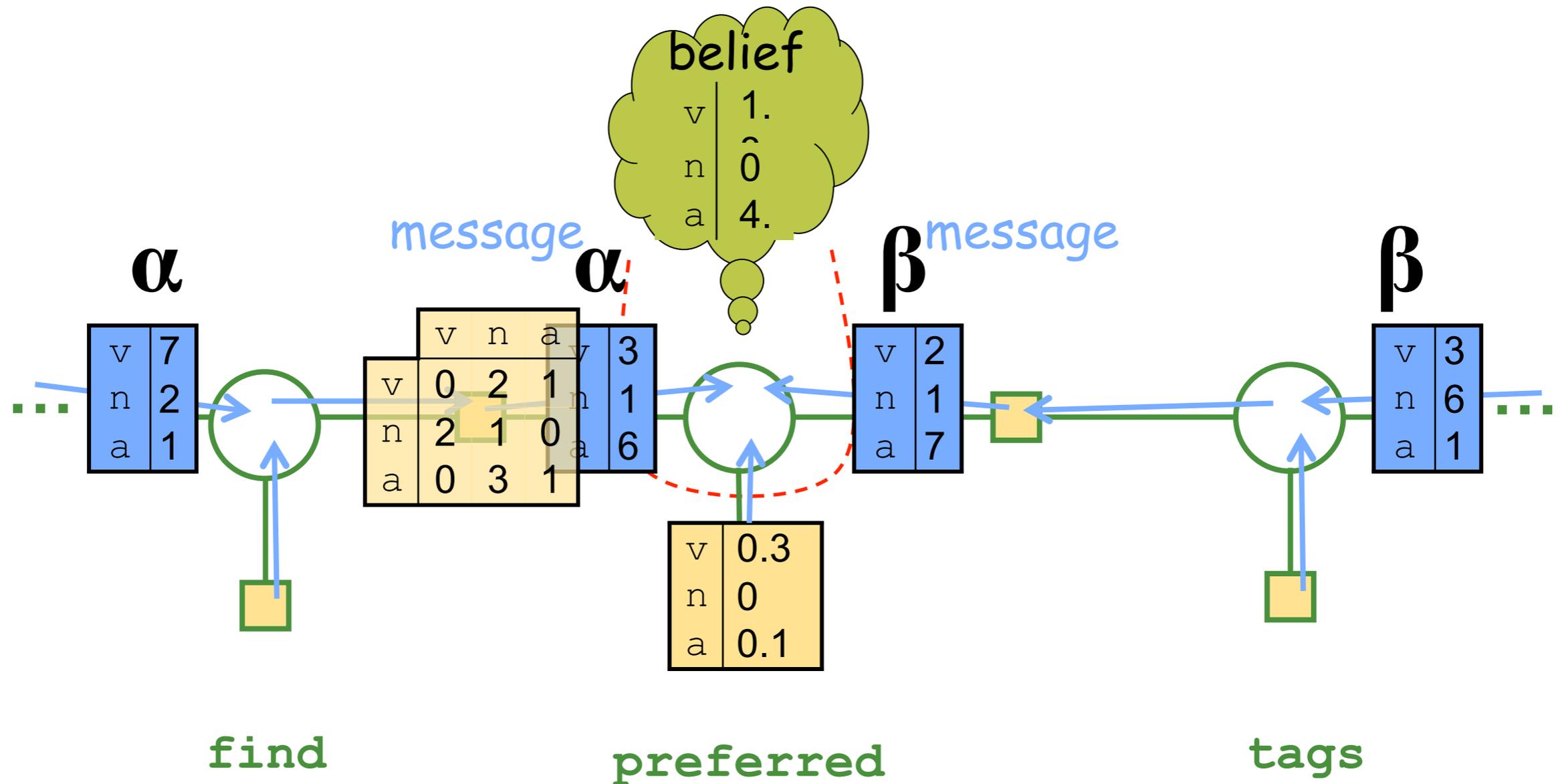
# Great ideas in ML: Forward-Backward

- In the CRF, message passing = forward-backward



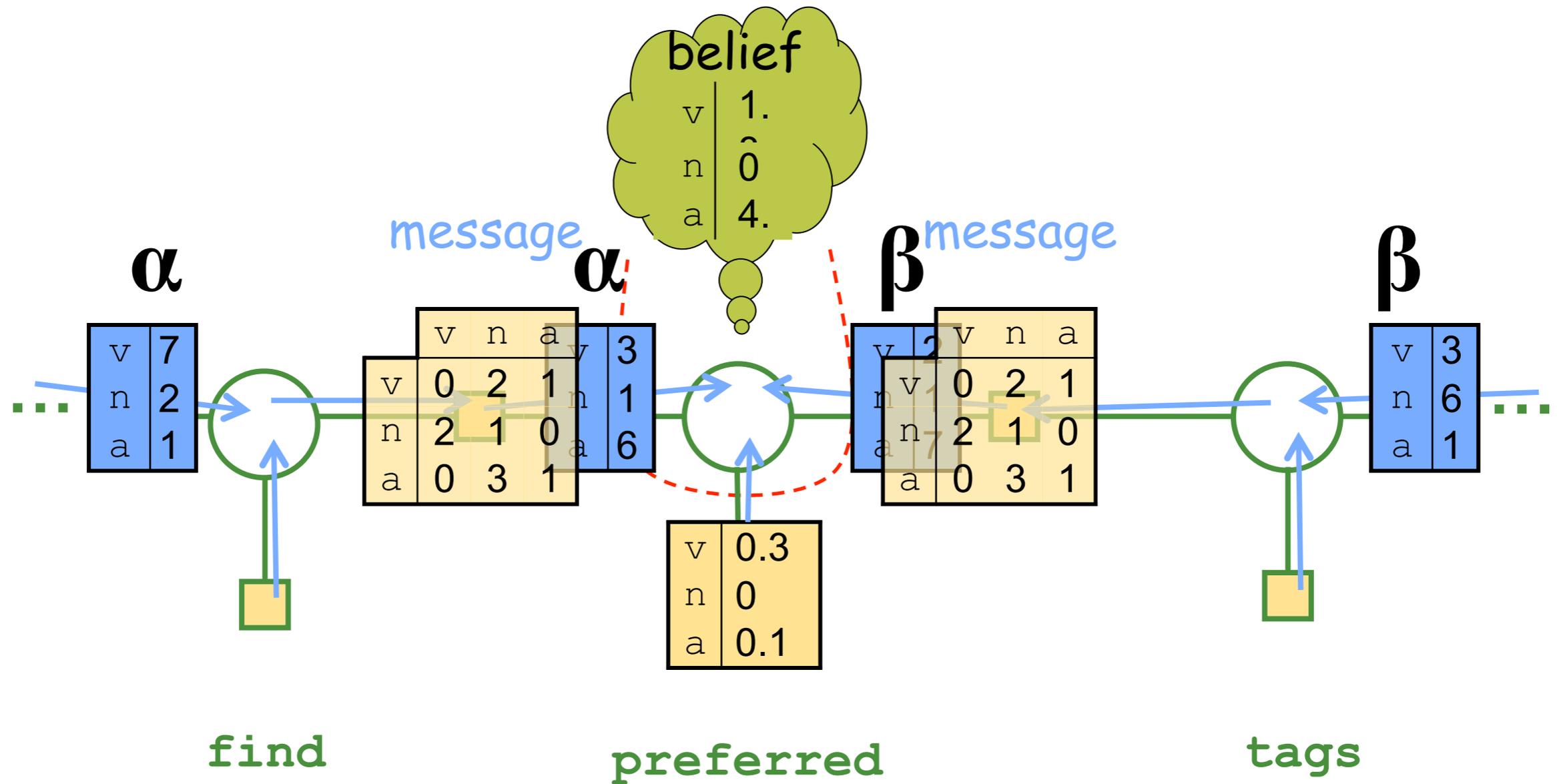
# Great ideas in ML: Forward-Backward

- In the CRF, message passing = forward-backward



# Great ideas in ML: Forward-Backward

- In the CRF, message passing = forward-backward

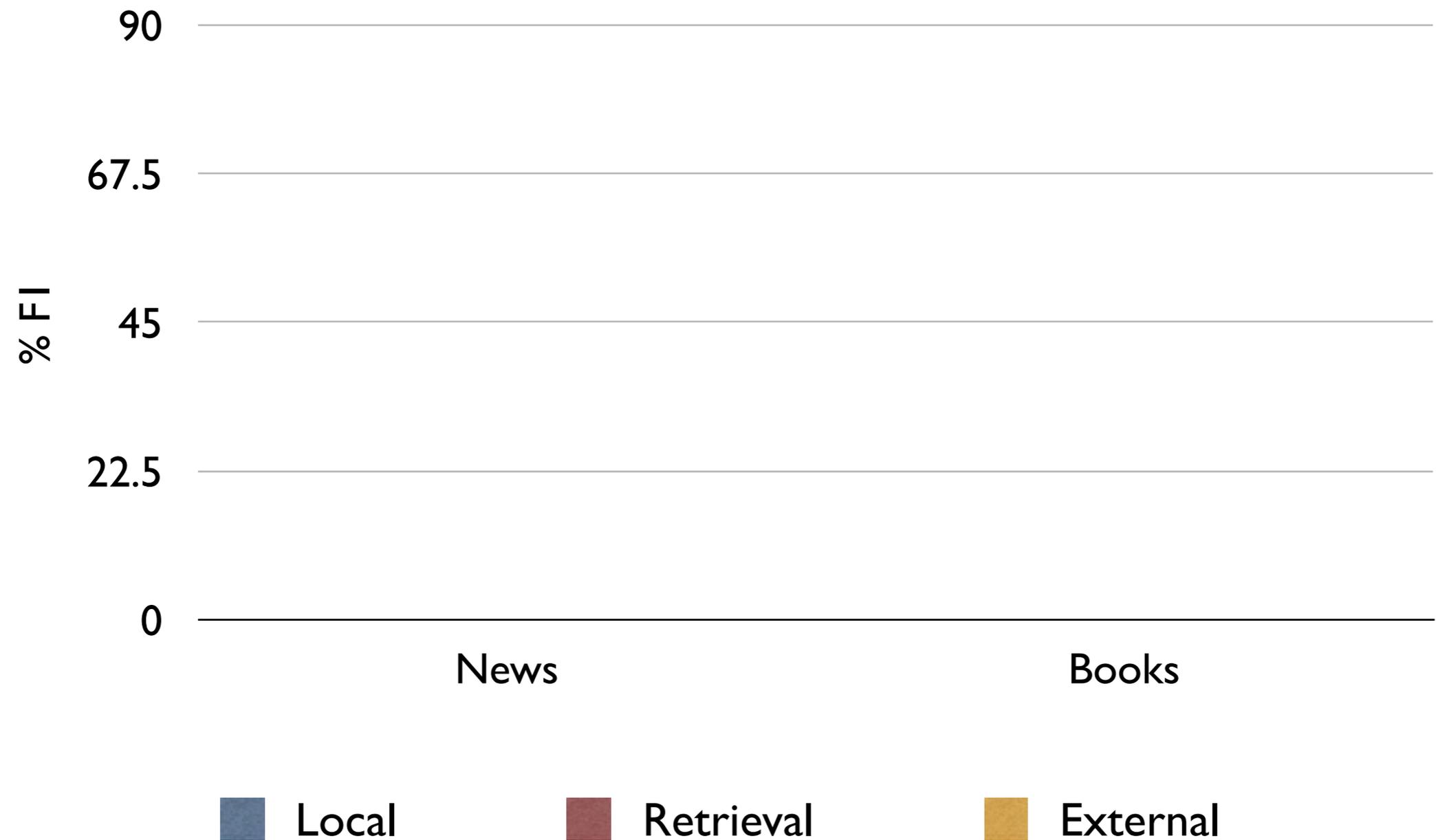


# Named Entity Recognition

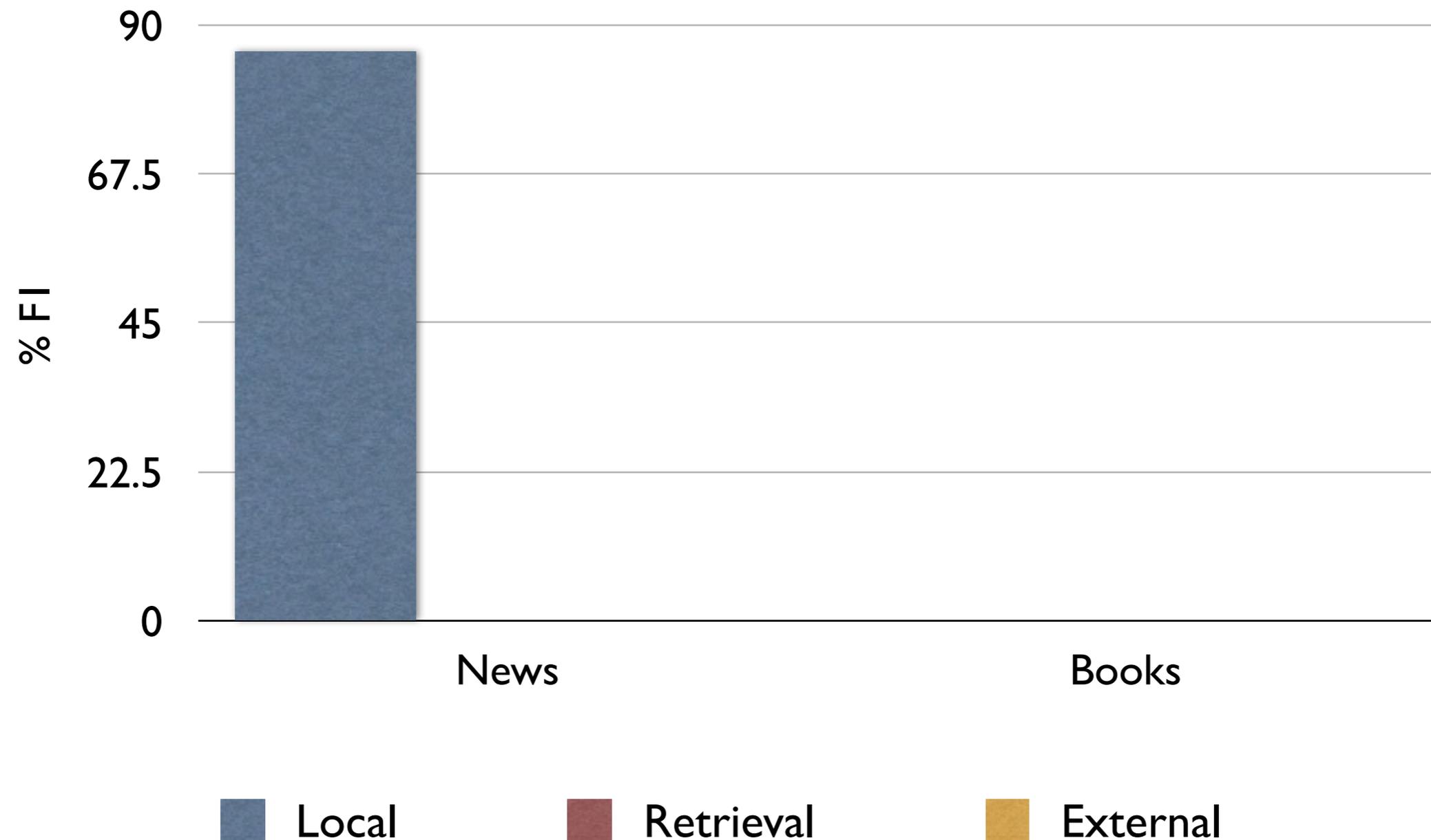
- Accurate recognition requires about 1M words of training data (1,500 news stories)
  - may be more expensive than developing rules for some applications
- Both rule-based and statistical can achieve about 90% effectiveness for categories such as names, locations, organizations
  - others, such as product name, can be much worse

# Domain Adaptation

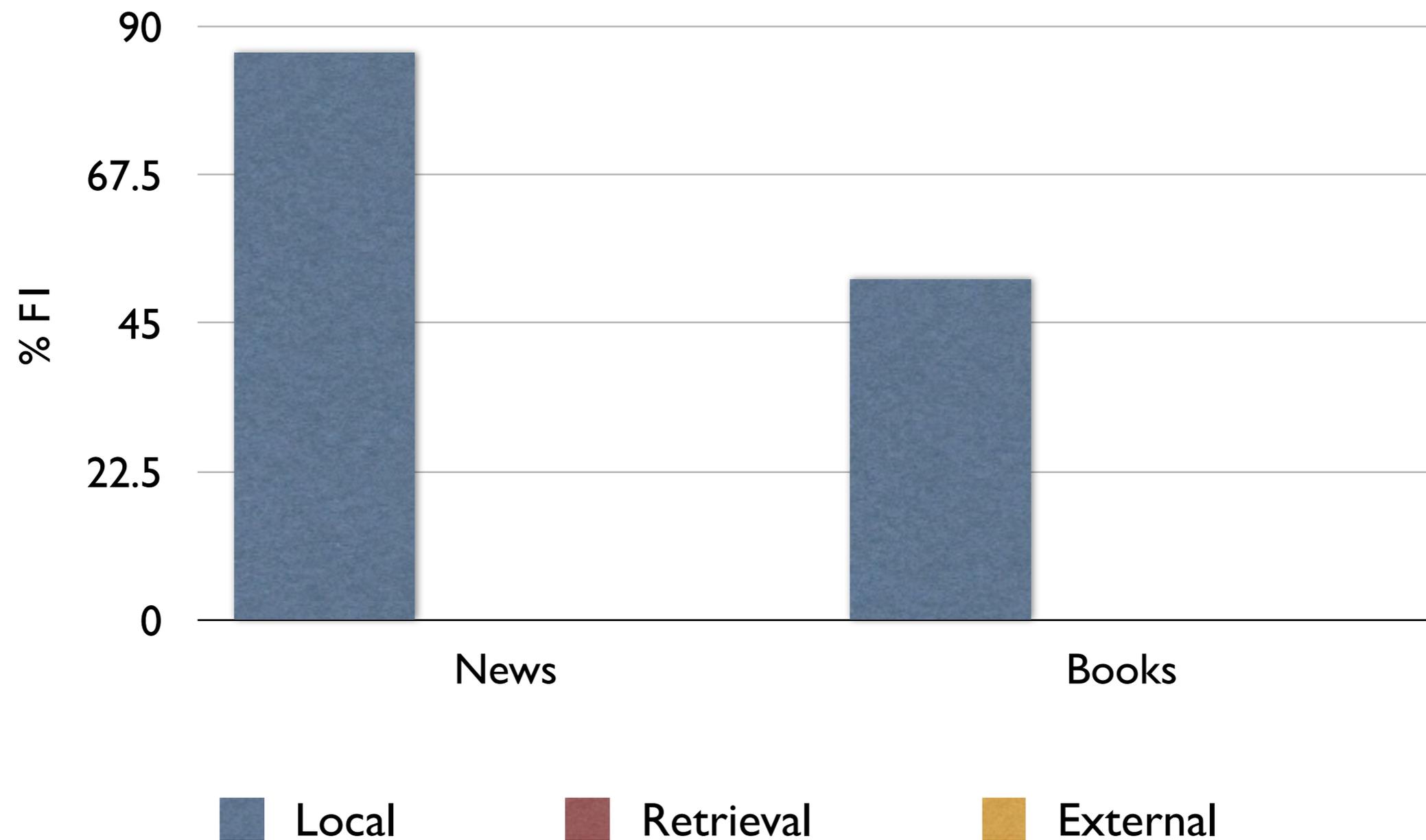
# Domain Adaptation



# Domain Adaptation



# Domain Adaptation

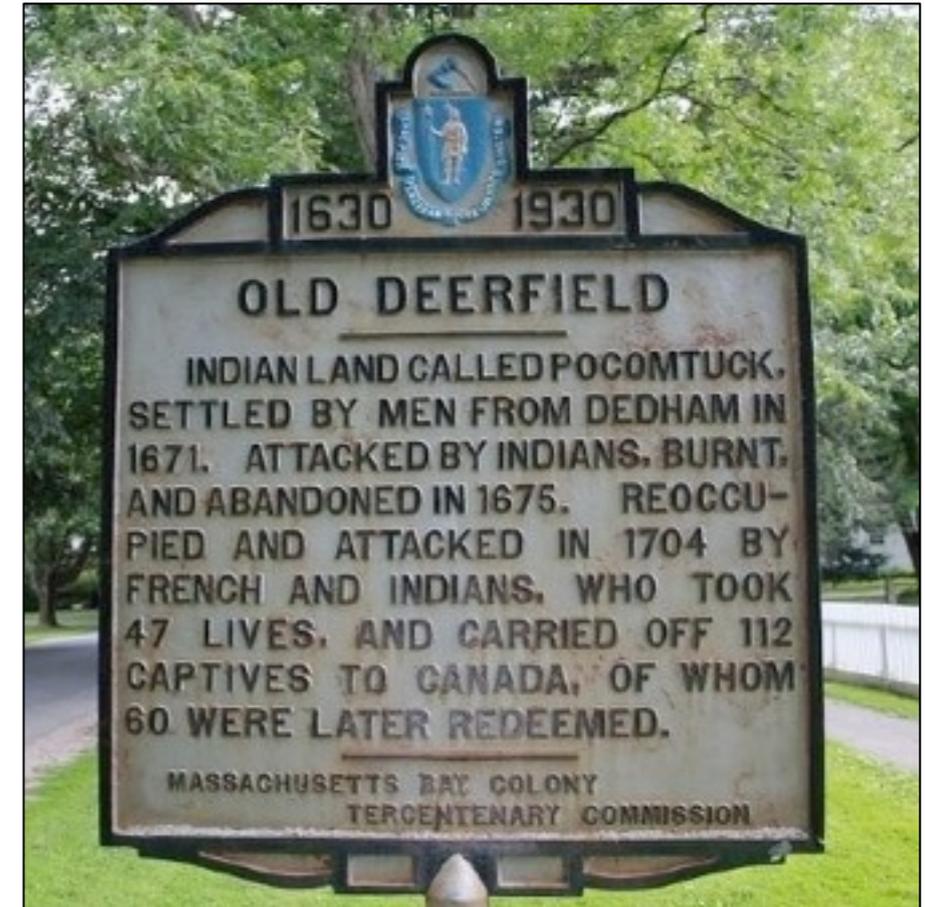


# Topic-specific NER

- Old Deerfield topic collection
  - 10 relevant books
  - 20 pages
  - Manually annotated entities

	<b>Count</b>
Tokens	10,050
Person	273
Miscellaneous	98
Location	241
Organization	49

Table 2: Historic Book NER Collection Statistics



# Errors in Book Data

- OCR Errors
  - NewEngland, District of Mains, Mudd } Brook, ye Enghsh  
Sugar loafe Hill
- Differences in style, caps, and punctuation
  - Connecticut river, Willard house, Liberty to hunt Deere  
or other Wild creatures and to gather Walnuts ,  
Chestnuts and other nuts things etc . on ye commons
- Sparsity in the features in the test tagging
  - Hatfield, Mount Tom and Mount Holyoke are all tagged as  
people!
- Inconsistent tagging
  - Nipmucks is tagged inconsistently in the same sentence!

# Passage Retrieval for Sequence Labeling

CIKM 2011

$X_i$ : [Distillery]

The Arran **Distillery** is a patron of the World Burns Federation and as such has created a Robert Burns Single Malt and Robert Burns Blended Whisky in honour of Scotland's National Poet

# Passage Retrieval for Sequence Labeling

CIKM 2011

$X_i$ : [Distillery]

## Query Q:

```
#weight( 0.8 #combine (Arran Distillery)
0.15 #combine( #ow1(Arran Distillery)) )
0.05 #combine(#ow8(Arran Distillery)) )
```

The Arran **Distillery** is a patron of the World Burns Federation and as such has created a Robert Burns Single Malt and Robert Burns Blended Whisky in honour of Scotland's National Poet

# Passage Retrieval for Sequence Labeling

CIKM 2011

R: Ranked list of similar sequences

$X_i$ : [Distillery]

Query Q:

```
#weight( 0.8 #combine (Arran Distillery)
0.15 #combine( #ow1(Arran Distillery)) )
0.05 #combine(#ow8(Arran Distillery)) )
```

0.5

Arran **Distillery** was founded in 1995 by Harold Currie, a former Managing Director at Chivas

0.2

Arran Single Malt is a Single Malt Scotch whisky distilled by the Arran **Distillery**, the only

0.2

Picturesque Lochranza, at the north of the Island is the location of Arran's first legal

0.17

The Isle of Arran **Distillery** produces a light, aromatic single malt.

0.15

Isle of Arran **distillery**, the birth-place of the award-winning Arran Single Malt whisky.

0.07

There are hairy coos at the Isle of Arran **distillery**.

# Passage Retrieval for Sequence Labeling

CIKM 2011

**R:** Ranked list of similar sequences

$X_i$ : [Distillery]

**Query Q:**

```
#weight( 0.8 #combine (Arran Distillery)
0.15 #combine( #ow1(Arran Distillery)) )
0.05 #combine(#ow8(Arran Distillery)) )
```



0.5

Arran **Distillery** was founded in 1995 by Harold Currie, a former Managing Director at Chivas

0.2

Arran Single Malt is a Single Malt Scotch whisky distilled by the Arran **Distillery**, the only

0.2

Picturesque Lochranza, at the north of the Island is the location of Arran's first legal

0.17

The Isle of Arran **Distillery** produces a light, aromatic single malt.

0.15

Isle of Arran **distillery**, the birth-place of the award-winning Arran Single Malt whisky.

0.07

There are hairy coos at the Isle of Arran **distillery**.

# Using Evidence from Passages

CIKM 2011

- Weighted feature copying to estimate  $p(y_i | x_i')$ 
  - Aggregation of feature probabilities in retrieved set

$P(\theta_r | Q_{x_i})$

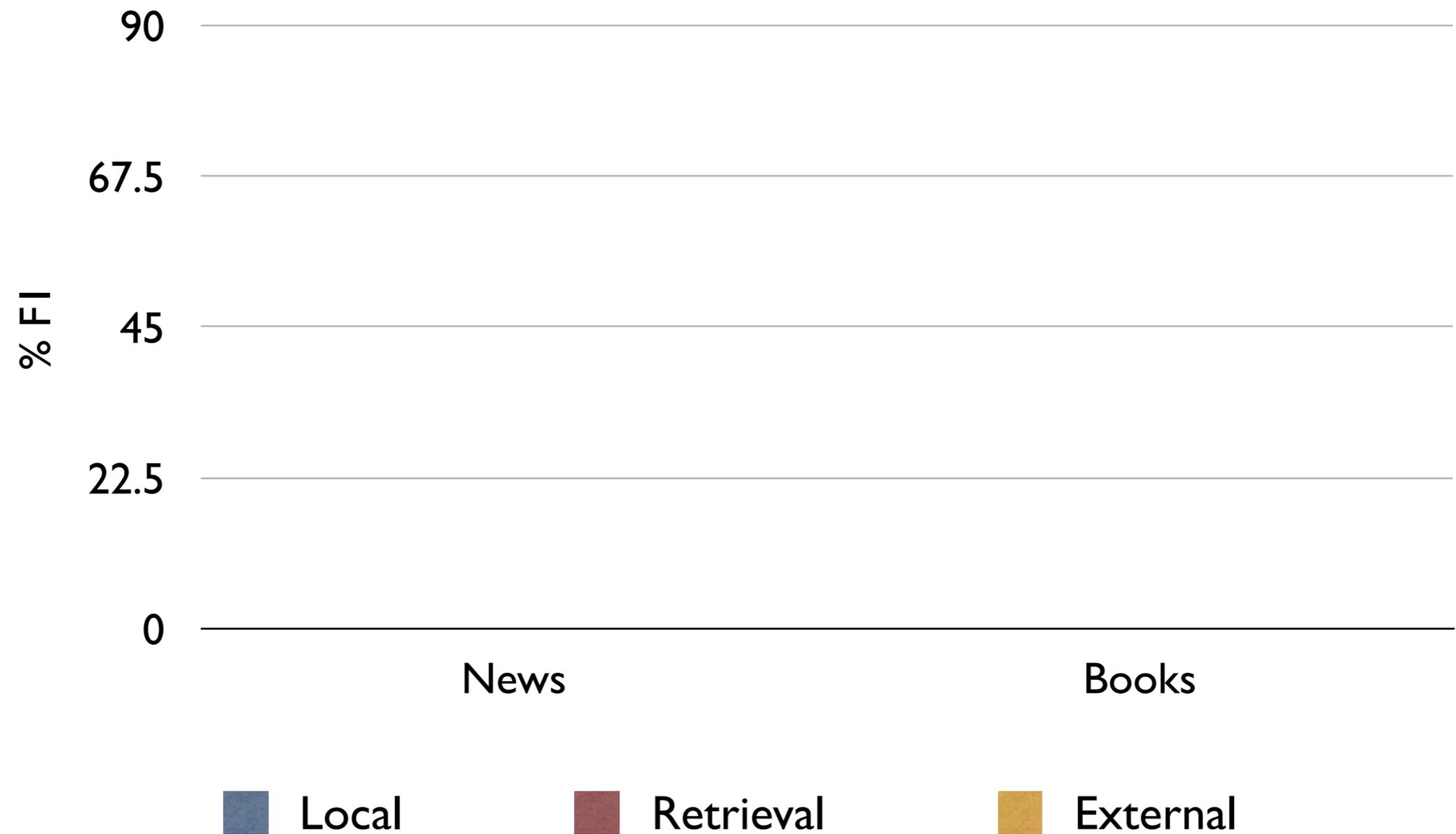
0.5	Arran <b>Distillery</b> was founded in 1995 by Harold Currie, a former Managing
0.2	Arran Single Malt is a Single Malt Scotch whisky distilled by the Arran <b>Distillery</b> ,
0.2	Picturesque Lochranza, at the north of the Island is the location of Arran's first legal <b>distillery</b> for over 150 years.
0.17	The Isle of Arran <b>Distillery</b> produces a light, aromatic single malt.
0.15	Isle of Arran <b>distillery</b> , the birth-place of the award-winning Arran Single Malt
0.07	There are hairy coos at the Isle of Arran <b>distillery</b> .
0.02	The Macallan <b>Distillery</b> is located in Craigellachie in the Speyside ..

Aggregate Features	Val
PREV_WORD_arran	1.0
RET_PREV_WORD_arran	0.7320
RET_PREV_WORD_macallan	0.015
RET_PREV_WORD_legal	0.15
CUR_CAP_CAPITALIZED	1.0
RET_CUR_CAP_CAPITALIZED	0.83

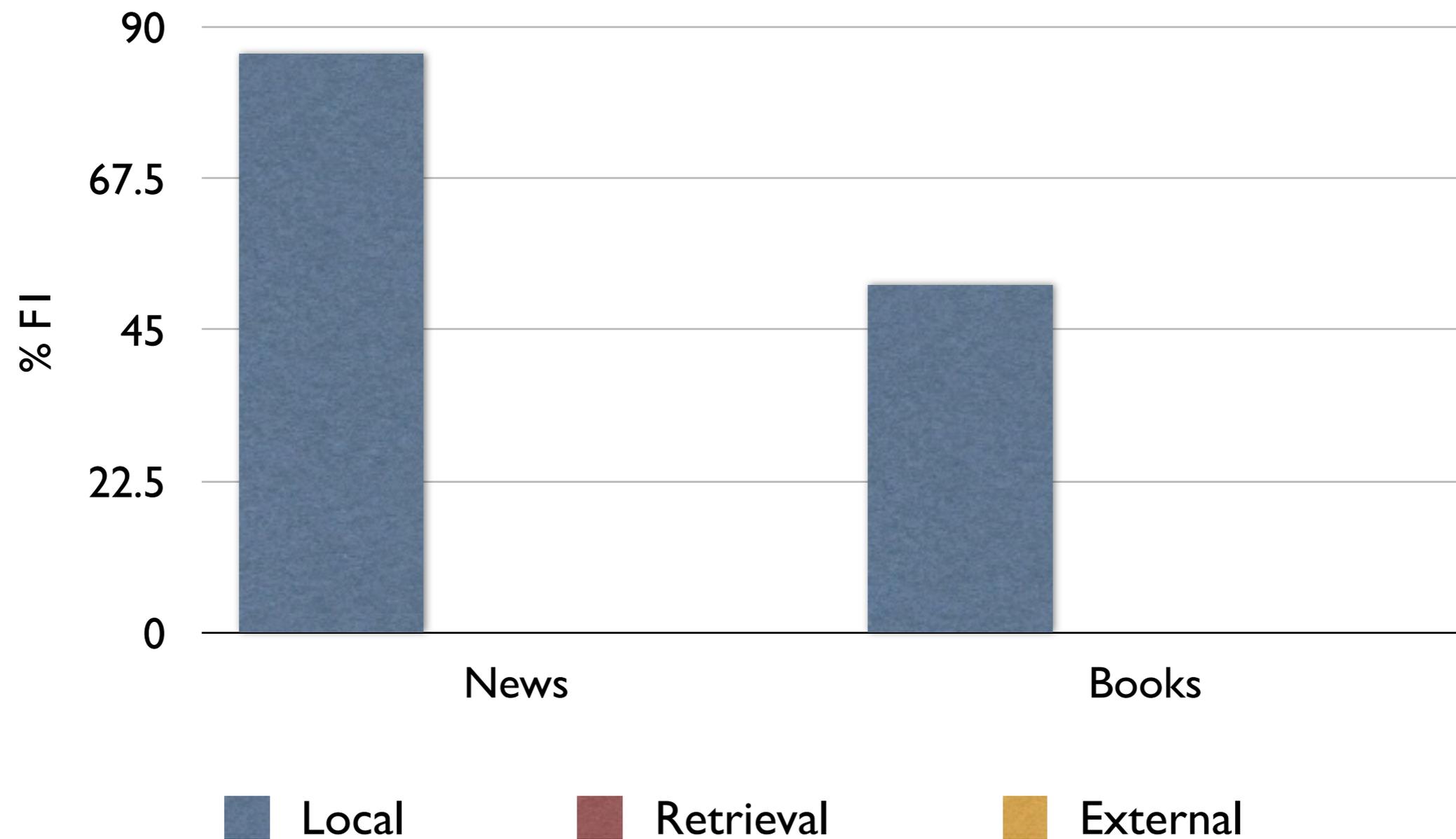
*Separate features from retrieved passages (RET)*

# Domain Adaptation

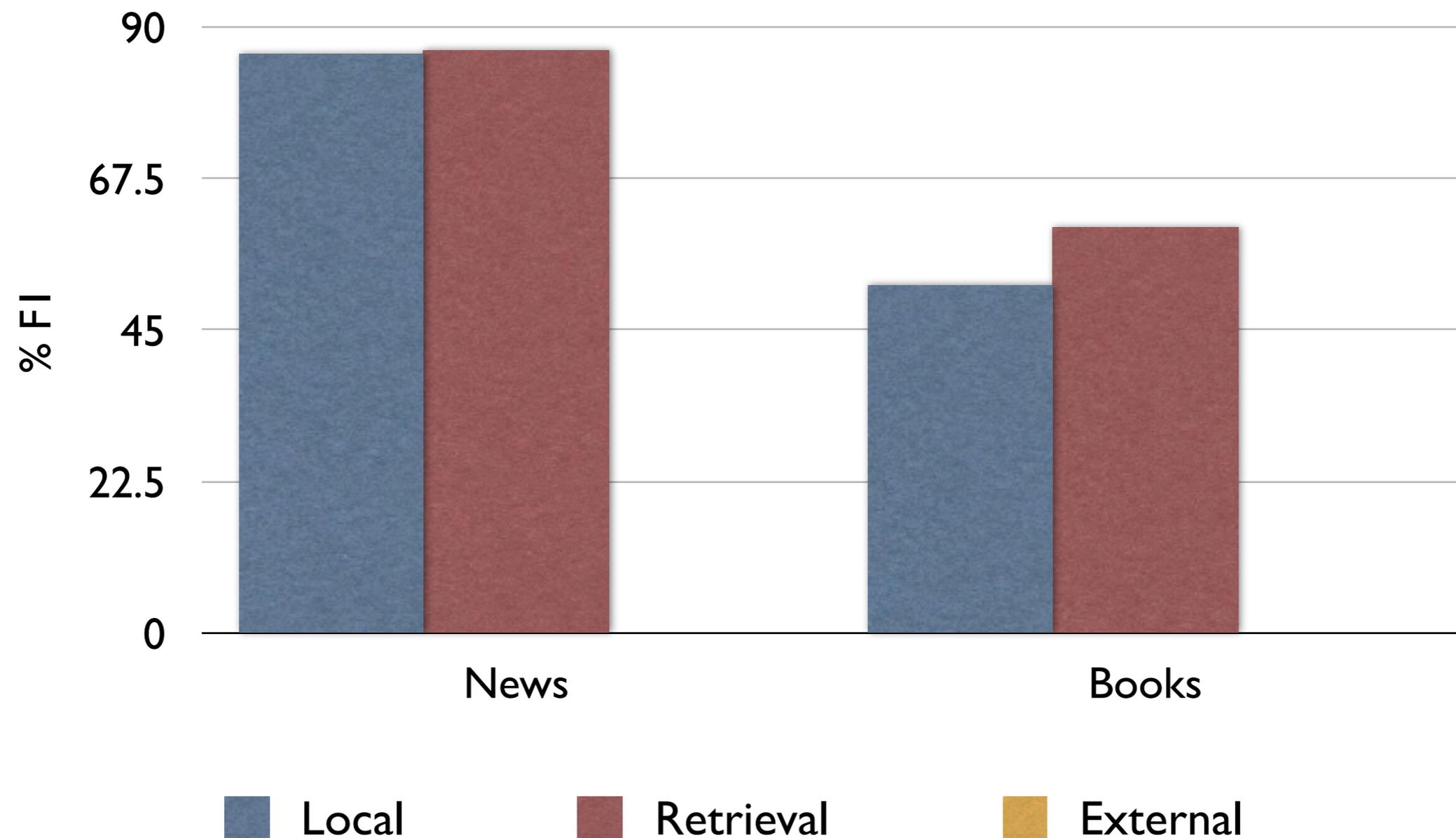
# Domain Adaptation



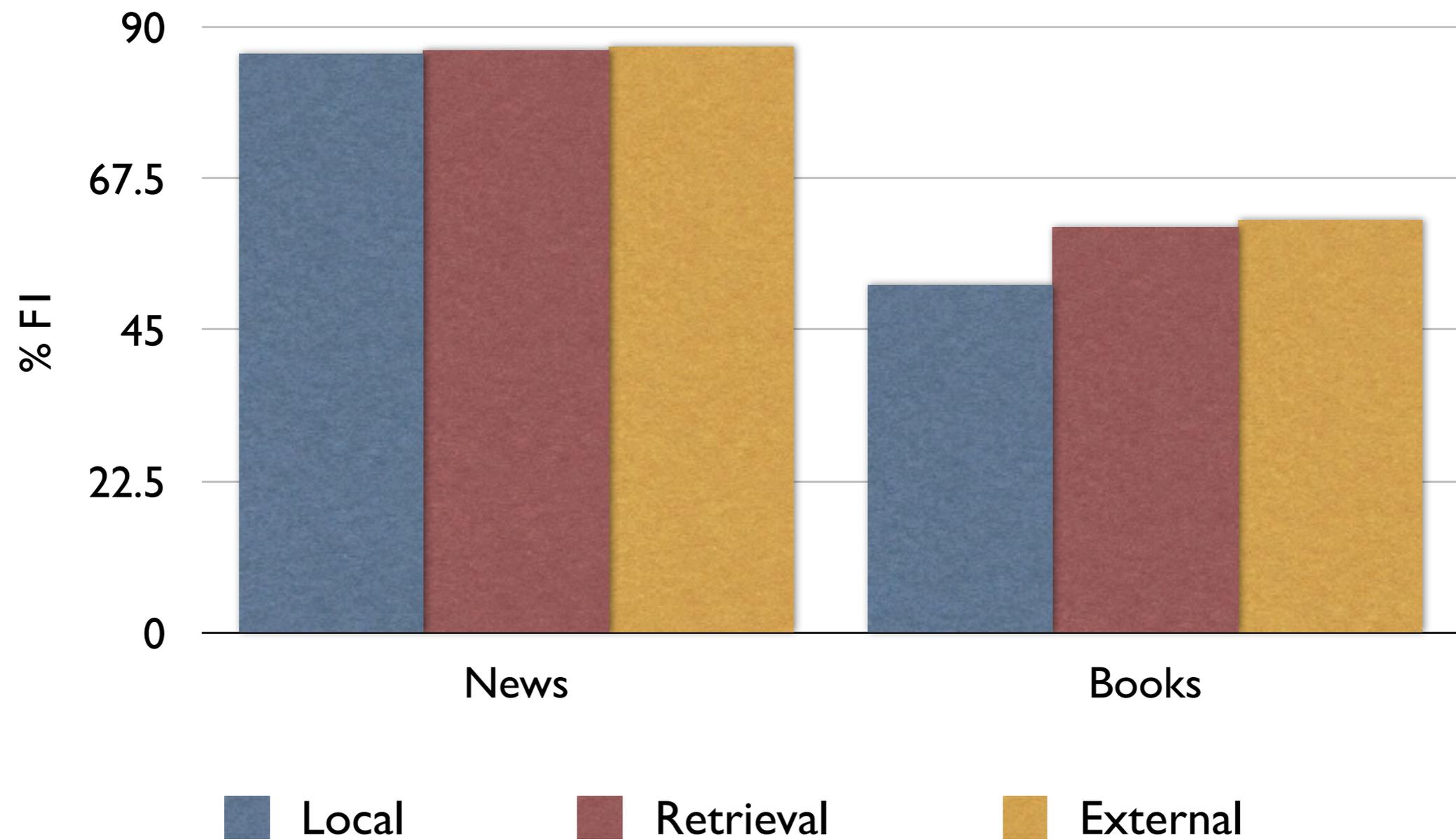
# Domain Adaptation



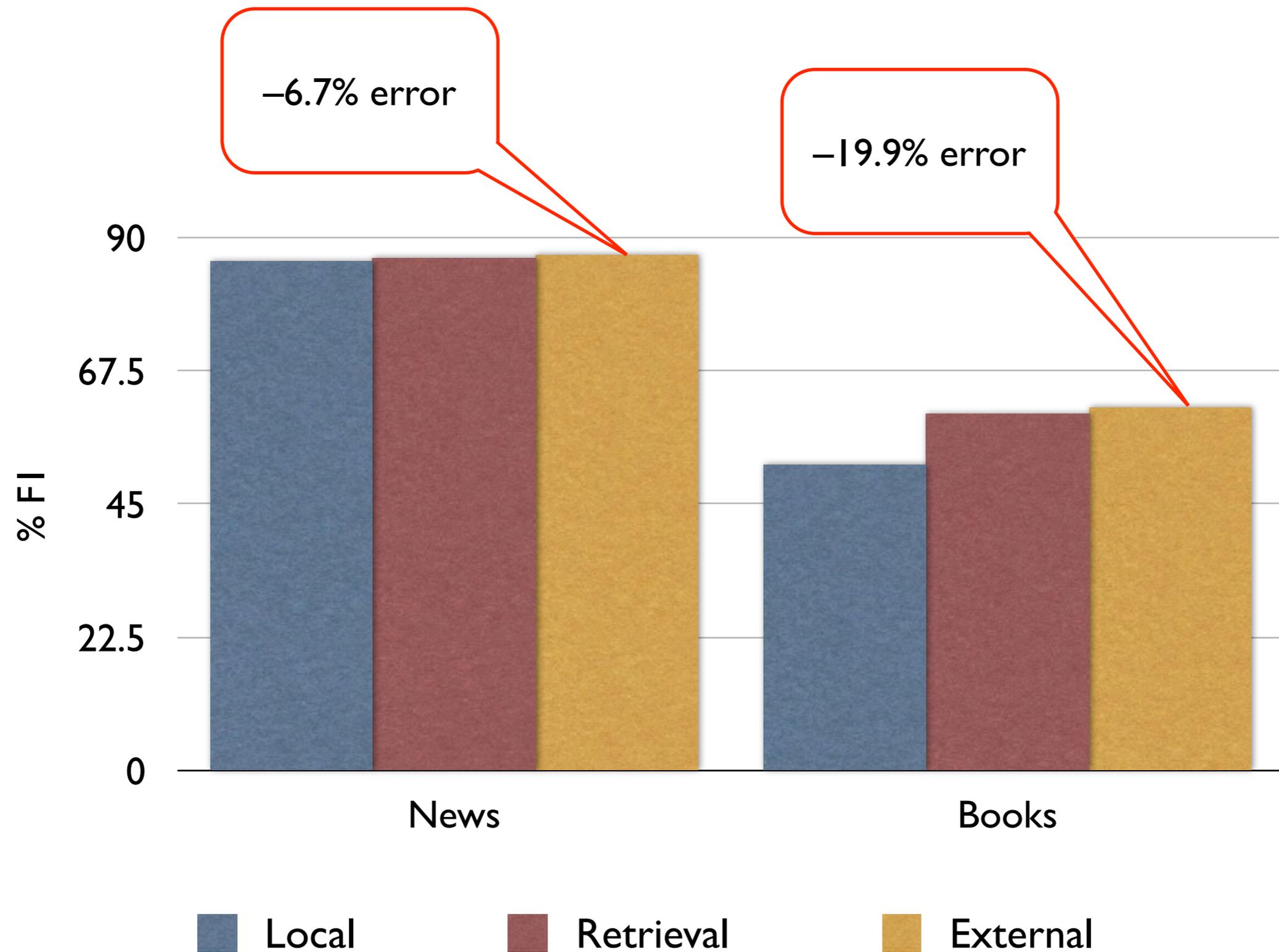
# Domain Adaptation



# Domain Adaptation



# Domain Adaptation



# Internationalization

- 2/3 of the Web is in English
- About 50% of Web users do not use English as their primary language
- Many (maybe most) search applications have to deal with multiple languages
  - monolingual search*: search in one language, but with many possible languages
  - cross-language search*: search in multiple languages at the same time

# Internationalization

- Many aspects of search engines are language-neutral
- Major differences:
  - Text encoding (converting to Unicode)
  - Tokenizing (many languages have no word separators)
  - Stemming
- Cultural differences may also impact interface design and features provided

# Chinese “Tokenizing”

## 1. Original text

旱灾在中国造成的影响

(the impact of droughts in China)

## 2. Word segmentation

旱灾 在 中国 造成 的 影响

drought at china make impact

## 3. Bigrams

旱灾 灾在 在中 中国 国造

造成 成的 的影 影响