



Chapter 5 Natural Language Processing

COMP 6721 Introduction of AI

Russell & Norvig – Chapter 23.1 + 23.2 + 23.3

Topics

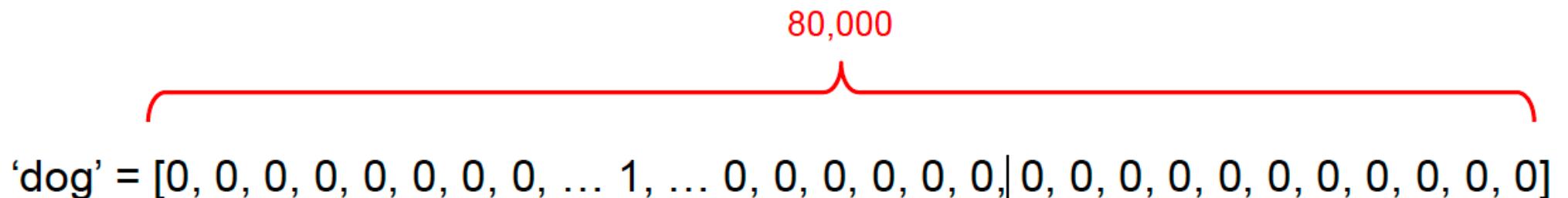
- ▶ Word Embeddings
- ▶ N-gram Models
- ▶ Word Embeddings (continued)
- ▶ Contextualized Word Embeddings

Beyond One Hot Vectors

- ▶ One hot vectors are local representations.
- ▶ They represent a 1:1 mapping between a word (e.g., 'cat', 'dog', 'house') and a position in a vector (i.e., the position with the value of 1).

Beyond One Hot Vectors

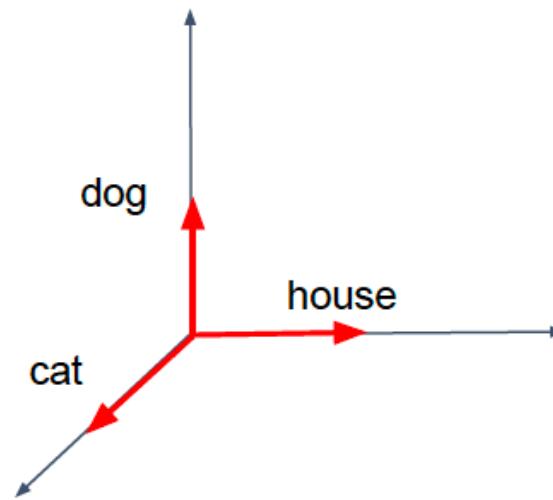
- ▶ This representation has two main disadvantages:
 - ▶ The vector dimension is proportional to the vocabulary size.
 - ▶ It does not capture relationships between words.
- ▶ E.g., with a vocabulary containing 80,000 words:



Beyond One Hot Vectors

- ▶ This representation has two main disadvantages:
 - ▶ The vector dimension is proportional to the vocabulary size.
 - ▶ It does not capture relationships between words.
 - ▶ All words are at the same distance.

$$\begin{array}{ll} \text{'cat'} & = [0, 0, 1] \\ \text{'dog'} & = [0, 1, 0] \\ \text{'house'} & = [1, 0, 0] \end{array}$$



Distributed Representation

- ▶ A different approach is to distribute the representation of a word across all available dimensions in a continuous space in a way that leads to “similar” words being close to each other.
- ▶ This is called a distributed representation.

Distributed Representation

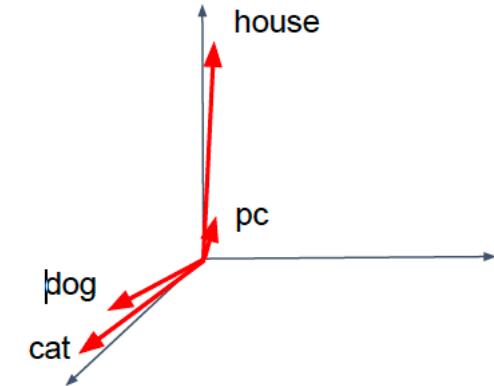
- Let's assume that the space dimensionality is 3 and that we have the following distributed representations:

'cat' = [1.1, 0.46, 45]

'dog' = [1.1, 1.10, 30]

'house' = [0.2, 4.51, 0]

'pc' = [0.1, 0.3 , 0]

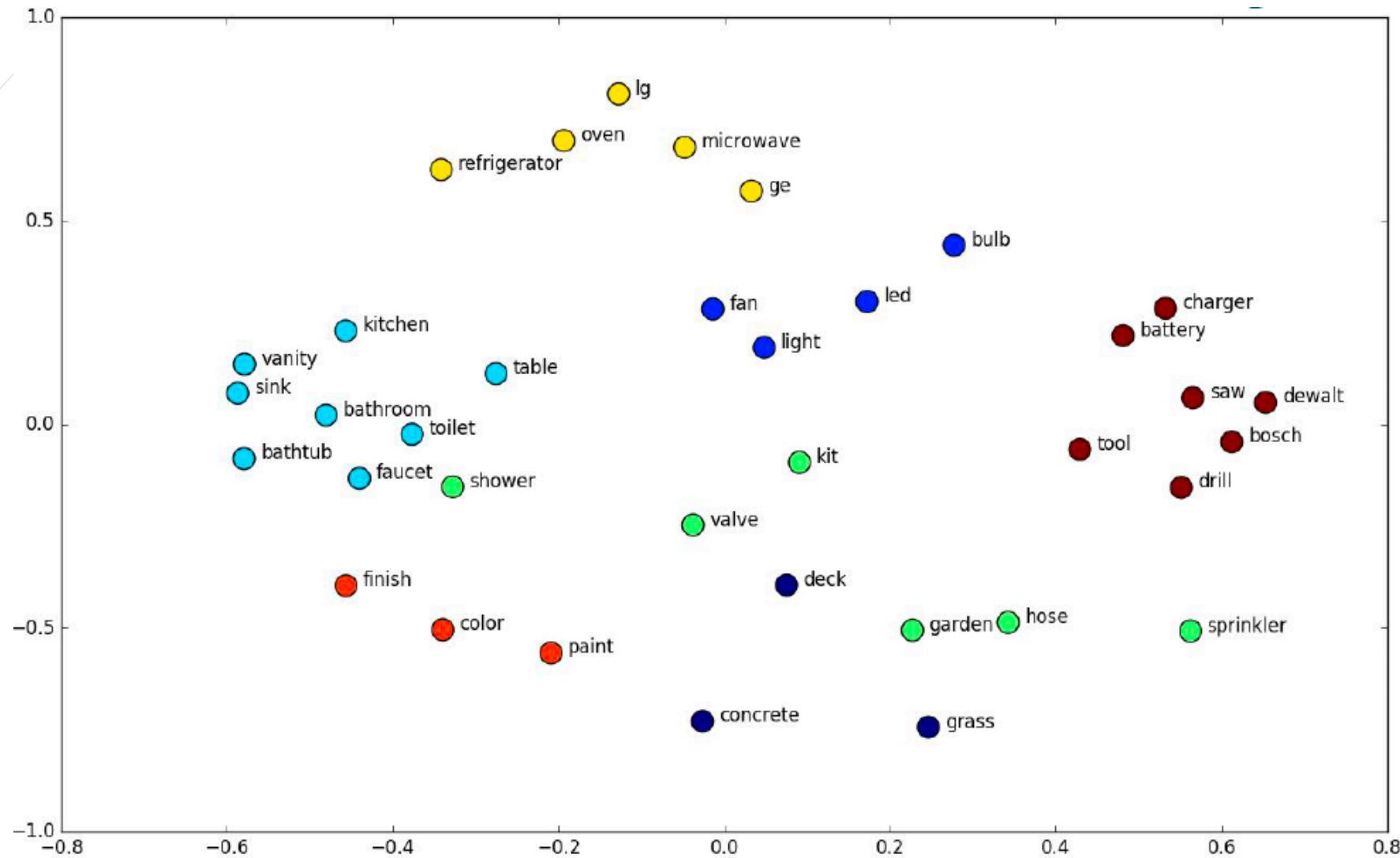


- Note how the vector size is not equal to the vocabulary size.
 - The vector size will in general be much smaller than the vocabulary size.
 - Also note how 'dog' and 'cat' are more "similar" than 'dog' and 'house'.

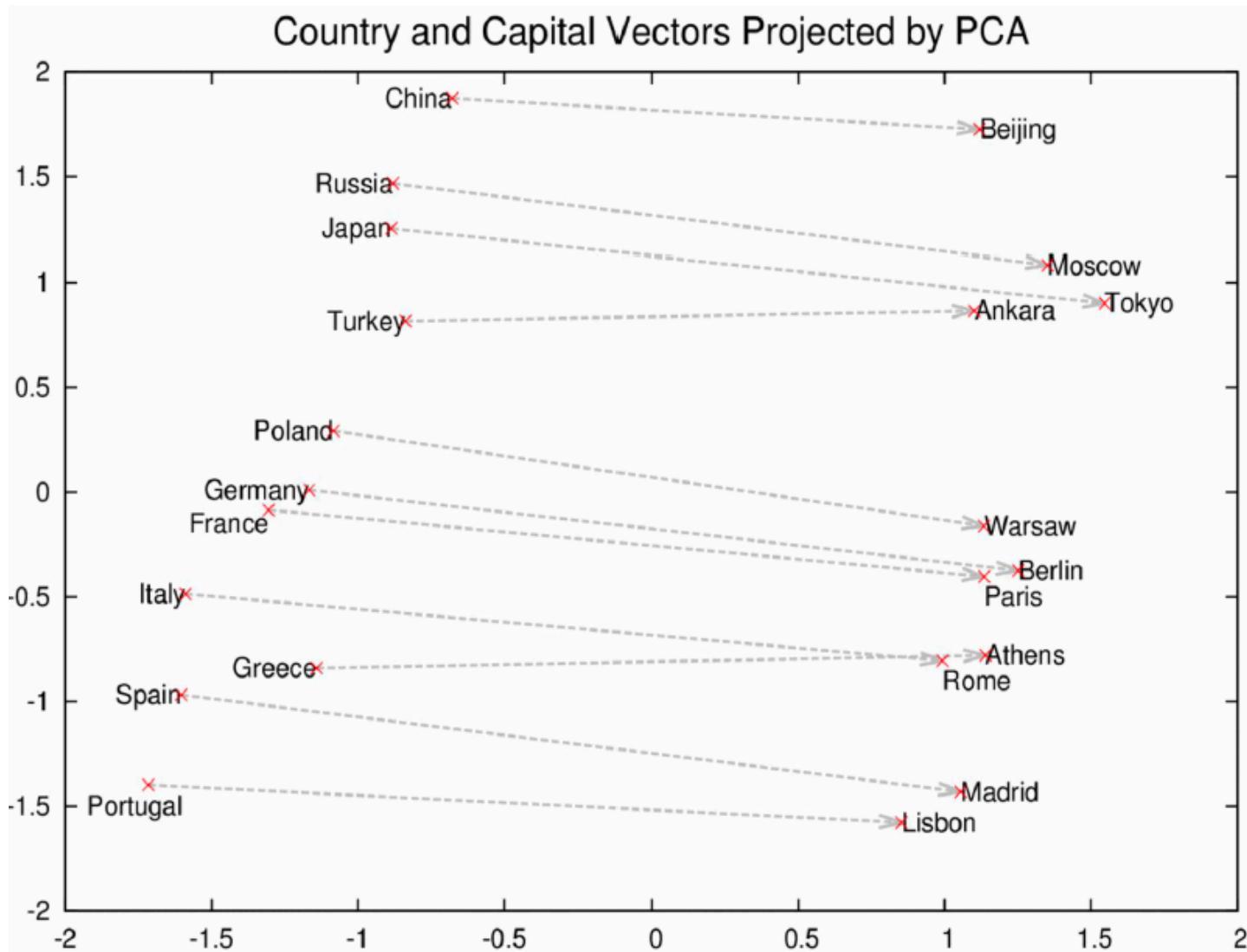
Distributed Representation

- ▶ Distributed representations can expose useful properties:
 - ▶ Semantically similar words can be close to each other.
 - ▶ Relationships between words can be captured.

Distributed Representation - Similarity



Distributed Representation - Relations



Word Embeddings

- ▶ Distributed representations for words are called word embeddings.
- ▶ There are several ways to generate word embeddings that expose the desired properties that we mentioned, including:
 - ▶ Word2Vec
 - ▶ FastText

Topics

- ▶ ***Word Embeddings***
- ▶ N-gram Models
- ▶ Word Embeddings (continued)
- ▶ Contextualized Word Embeddings

N-gram Model

- ▶ An n-gram model is a probability distribution over sequences of events (grams/units/items)
- ▶ models the order of the events
- ▶ Used when the past sequence of events is a good indicator of the next event to occur in the sequence
- ▶ i.e. To predict the next event in a sequence of event

N-gram Model

► E.g.

- next move of player based on his/her past moves
left right right up ... up? down? left? right?
- next base pair based on past DNA sequence
AGCTTCG ... A? G? C? T?
- next word based on past words
Hi dear, how are ... helicopter? laptop? you? magic?

What is a Language Model

- ▶ A Language model is a n-gram model over word/character sequences
- ▶ ie: events = words or events = character
- ▶ $P(\text{"I'd like a coffee with 2 sugars and milk"}) \approx 0.001$
- ▶ $P(\text{"I'd hike a toffee with 2 sugars and silk"}) \approx 0.000000001$

Applications of Language Model

- ▶ Speech Recognition
- ▶ Statistical Machine Translation
- ▶ Language Identification
- ▶ Spelling correction

He is trying to fine out.
He is trying to find out.
- ▶ Optical character recognition / Handwriting recognition
- ▶ ...

In Speech Recognition



Given: Observed sound - O

Find: The most likely word/sentence - S^*

S_1 : How to recognize speech. ?

S_2 : How to wreck a nice beach. ?

S_3 : ...

- Goal: find most likely sentence (S^*) given the observed sound (O) ...

- i.e. pick the sentence with the highest probability:

$$S^* = \underset{S \in L}{\operatorname{argmax}} P(S | O)$$

- We can use Bayes rule to rewrite this as:

$$S^* = \underset{S \in L}{\operatorname{argmax}} \frac{P(O | S)P(S)}{P(O)}$$

- Since denominator is the same for each candidate S, we can ignore it for the argmax:

$$S^* = \underset{S \in L}{\operatorname{argmax}} P(O | S) \times P(S)$$

Acoustic model --

Probability of the possible phonemes in the language +

Probability of ≠ pronunciations

Language model -- $P(\text{a sentence})$

Probability of the candidate sentence in the language

In Speech Recognition

$$S^* = \underset{S \in L}{\operatorname{argmax}} P(O | S) \times P(S)$$

Acoustic model

Language model

argmax $P(\text{word sequence} | \text{acoustic signal})$

word sequence

$\operatorname{argmax}_{\text{wordsequence}}$
$$\frac{P(\text{acoustic signal} | \text{word sequence}) \times P(\text{word sequence})}{P(\text{acoustic signal})}$$

$\operatorname{argmax}_{\text{wordsequence}}$ $P(\text{acoustic signal} | \text{word sequence}) \times P(\text{word sequence})$

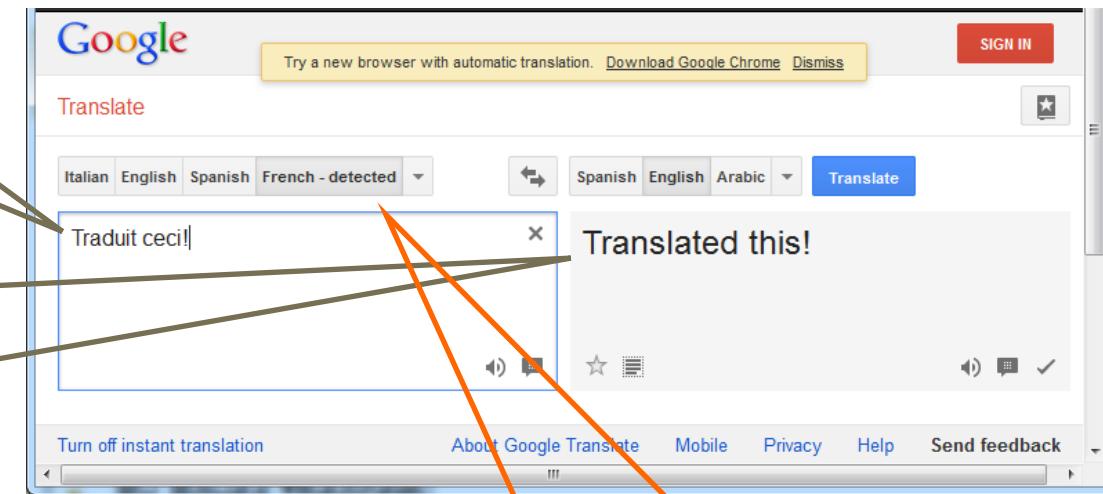
In statistical Machine Translation

→ Assume we translate from fr[foreign] to English i.e: (en | fr)

Given: Foreign sentence - fr

Find: The most likely English sentence - en*

- S1: Translate that!
- S2: Translate this!
- S3: Eat your soup!
- S4...



Automatic Language Identification...
guess how that's done?

$$\underset{\text{en}}{\text{en}^*} = \operatorname{argmax}_{\text{en}} P(\text{fr} | \text{en}) \times P(\text{en})$$

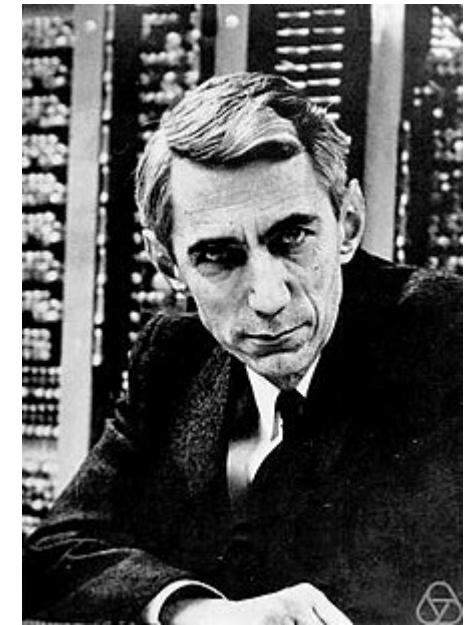
Translation model

Language model

“Shannon Game” (Shannon, 1951)

- ▶ Predict the next word/character given the $n-1$ previous words/characters.

“I am going to make a collect ...”



1st approximation

- ▶ each word has an equal probability to follow any other
 - ▶ with 100,000 words, the probability of each word at any given point is .00001
- ▶ but some words are more frequent than others...
 - ▶ “the” appears many more times, than “rabbit”

2nd approximation: unigrams

- ▶ take into account the frequency of the word in some training corpus
 - ▶ at any given point, “the” is more probable than “rabbit”
- ▶ but does not take word order into account. This is the **bag of word** approach.
 - ▶ *“Just then, the white ...”*
- ▶ so the probability of a word also depends on the previous words (the history)

$$P(w_n | w_1 w_2 \dots w_{n-1})$$

n-grams

- ▶ “*the large green _____.*”
 - ▶ “mountain”? “tree”?
- ▶ “*Sue swallowed the large green _____.*”
 - ▶ “pill”? “broccoli”?
- ▶ Knowing that Sue “swallowed” helps narrow down possibilities
- ▶ ie. Going back 3 words before helps
- ▶ But, how far back do we look?

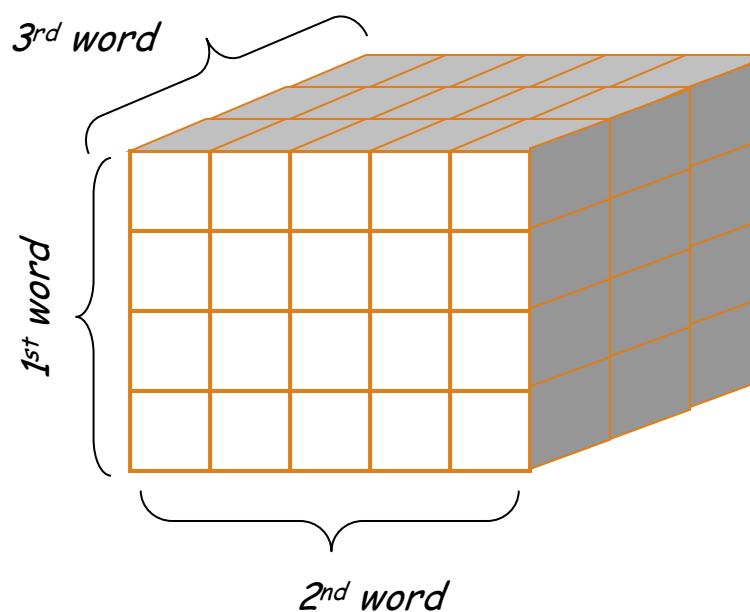
Bigrams

- ▶ first-order Markov models $P(w_n | w_{n-1})$
- ▶ N-by-N matrix of probabilities/frequencies
- ▶ N = size of the vocabulary we are using

	a	aardvark	aardwolf	aback	...	zoophyte	zucchini
a	0	0	0	0	...	8	5
aardvark	0	0	0	0	...	0	0
aardwolf	0	0	0	0	...	0	0
aback	26	1	6	0	...	12	2
...
zoophyte	0	0	0	1	...	0	0
zucchini	0	0	0	3	...	0	0

Trigrams

- ▶ second-order Markov models $P(w_n | w_{n-1} w_{n-2})$
- ▶ N-by-N-by-N matrix of probabilities/frequencies
- ▶ N = size of the vocabulary we are using



Why use only bi- or tri-grams

- ▶ Markov approximation is still costly with a 20 000 word vocabulary:
 - ▶ bigram needs to store 400 million parameters
 - ▶ trigram needs to store 8 trillion parameters
 - ▶ using a language model > trigram is impractical

Building N-gram Models

1. Data preparation:

- ▶ Decide on training corpus
- ▶ Clean and tokenize
- ▶ How do we deal with sentence boundaries?
 - ▶ I eat. I sleep.
 - ▶ (I eat) (eat I) (I sleep)
 - ▶ <s>I eat <s> I sleep <s>
 - ▶ (<s> I) (I eat) (eat <s>) (<s> I) (I sleep) (sleep <s>)

Building N-gram Models

2. Count words and build model
 - ▶ Let $C(w_1 \dots w_n)$ be the frequency of n-gram $w_1 \dots w_n$

$$P(w_n | w_1 \dots w_{n-1}) = \frac{C(w_1 \dots w_n)}{C(w_1 \dots w_{n-1})}$$

3. Smooth your model (see later)

Example 1:

- ▶ in a training corpus, we have 10 instances of “come across”
 - ▶ 8 times, followed by “as”
 - ▶ 1 time, followed by “more”
 - ▶ 1 time, followed by “a”
- ▶ so we have:
 - ▶ $P(\text{as} \mid \text{come across}) = \frac{C(\text{come across as})}{C(\text{come across})} = \frac{8}{10}$
 - ▶ $P(\text{more} \mid \text{come across}) = 0.1$
 - ▶ $P(\text{a} \mid \text{come across}) = 0.1$
 - ▶ $P(X \mid \text{come across}) = 0$ where $X \neq \text{"as", "more", "a"}$

Example 2:

$P(\text{on} \text{eat}) = .16$	$P(\text{want} I) = .32$	$P(\text{eat} to) = .26$
$P(\text{some} \text{eat}) = .06$	$P(\text{would} I) = .29$	$P(\text{have} to) = .14$
$P(\text{British} \text{eat}) = .001$	$P(\text{don't} I) = .08$	$P(\text{spend} to) = .09$
...
$P(I <s>) = .25$	$P(to want) = .65$	$P(\text{food} \text{British}) = .6$
$P(I'd <s>) = .06$	$P(a want) = .5$	$P(\text{restaurant} \text{British}) = .15$
...

$P(I \text{ want to eat British food})$

$$\begin{aligned}
 &= P(I | <s>) \times P(want | I) \times P(to | want) \times P(eat | to) \times P(British | eat) \times P(food | British) \\
 &= .25 \quad \times .32 \quad \times .65 \quad \times .26 \quad \times .001 \quad \times .6 \\
 &= .000008
 \end{aligned}$$

Remember this slide...

Using logs

- ▶ if we really do the product of probabilities...
 - ▶ $\operatorname{argmax}_{c_j} P(c_j) \prod P(w_i | c_j)$
 - ▶ we soon have numerical underflow...
 - ▶ ex: $0.01 \times 0.02 \times 0.05 \times \dots$
- ▶ so instead, we add the log of the probs
 - ▶ $\operatorname{argmax}_{c_j} \log(P(c_j)) + \sum \log(P(w_i | c_j))$
 - ▶ ex: $\log(0.01) + \log(0.02) + \log(0.05) + \dots$

Remember this slide...

Example

► Dataset



c1: COOKING

doc₁: ... stove... kitchen... the... heat
doc₂: ... kitchen... pasta... stove...
...
doc₁₀₀₀₀₀: ... stove...heat... ball...

c2: SPORTS

doc₁: ... ball... heat...
doc₂: ... the... referee... player...
...
doc₇₅₀₀₀: goal... injury ...

■ Assume:

- |V| = 100 vocabulary = {ball, heat, kitchen, referee, stove, the, ... }
- 500,000 words in Cooking
- 300,000 words in Sports
- 100,000 docs in Cooking
- 75,000 docs in Sports

Remember this slide...

Example

- ▶ Training – Unsmoothed / Smoothed probs:

▶ $P(\text{ball} \text{COOKING}) = \frac{10,000}{500,000}$	$\frac{??}{??}$	$P(\text{ball} \text{SPORTS}) = \frac{10,000}{300,000}$	$\frac{??}{??}$
▶ $P(\text{heat} \text{COOKING}) = \frac{255}{500,000}$	$\frac{??}{??}$	$P(\text{heat} \text{SPORTS}) = \frac{1,8000}{300,000}$	$\frac{??}{??}$
▶ $P(\text{kitchen} \text{COOKING}) = \frac{2,600}{500,000}$	$\frac{??}{??}$	$P(\text{kitchen} \text{SPORTS}) = \frac{0}{300,000}$	$\frac{??}{??}$
▶ $P(\text{referee} \text{COOKING}) = \frac{0}{500,000}$	$\frac{??}{??}$	$P(\text{referee} \text{SPORTS}) = \frac{1,500}{300,000}$	$\frac{??}{??}$
▶ $P(\text{stove} \text{COOKING}) = \frac{3,600}{500,000}$	$\frac{??}{??}$	$P(\text{stove} \text{SPORTS}) = \frac{4}{300,000}$	$\frac{??}{??}$
▶ $P(\text{the} \text{COOKING}) = \frac{400,000}{500,000}$	$\frac{??}{??}$	$P(\text{the} \text{SPORTS}) = \frac{19,000}{300,000}$	$\frac{??}{??}$
▶ ...			
▶ $P(\text{COOKING}) = \frac{100,000}{175,000}$		$P(\text{SPORTS}) = \frac{75,000}{175,000}$	

Remember this slide...

44

Example

- ▶ Training – Unsmoothed / Smoothed probs:
 - ▶ $P(COOKING) = \frac{100,000}{175,000}$ $P(SPORTS) = \frac{75,000}{175,000}$
- ▶ Testing: “the referee hit the ~~blue bird~~”
 - ▶ $\text{Score}(COOKING) = \log\left(\frac{100,000}{175,000}\right) + \log(P(\text{the} | COOKING)) + \log(P(\text{referee} | COOKING)) + \log(P(\text{hit} | COOKING)) + \log(P(\text{the} | COOKING))$
 - ▶ $\text{Score}(SPORTS) = \log\left(\frac{75,000}{175,000}\right) + \log(P(\text{the} | SPORTS)) + \log(P(\text{referee} | SPORTS)) + \log(P(\text{hit} | SPORTS)) + \log(P(\text{the} | SPORTS))$

Some Adjustments

- ▶ product of probabilities... numerical underflow for long sentences
- ▶ so instead of multiplying the probs, we add the log of the probs

$P(I \text{ want to eat British food})$

$$\begin{aligned} &= \log(P(I | <s>)) + \log(P(\text{want} | I)) + \log(P(\text{to} | \text{want})) + \log(P(\text{eat} | \text{to})) + \\ &\quad \log(P(\text{British} | \text{eat})) + \log(P(\text{food} | \text{British})) \\ &= \log(.25) + \log(.32) + \log(.65) + \log (.26) + \log(.001) + \log(.6) \end{aligned}$$

Problem: Data Sparseness

- ▶ What if a sequence never appears in training corpus? $P(X)=0$
 - ▶ “come across the men” --> prob = 0
 - ▶ “come across some men” --> prob = 0
 - ▶ “come across 3 men” --> prob = 0
- ▶ The model assigns a probability of zero to unseen events ...
- ▶ probability of an n-gram involving unseen words will be zero!
- ▶ Solution: smoothing
 - ▶ decrease the probability of previously seen events
 - ▶ so that there is a little bit of probability mass left over for previously unseen events

Remember this slide ...

40

Smooth Probabilities

- ▶ normally: $P(w_i | c_j) = \frac{\text{(frequency of } w_i \text{ in } c_j)}{\text{total number of words in } c_j}$
- ▶ what if we have a $P(w_i | c_j) = 0 \dots ?$
 - ▶ ex. the word "fake" never appeared in the class SPAM?
 - ▶ then $P(\text{"fake"} | \text{SPAM}) = 0$
- ▶ so if a text contains the word "fake", the class SPAM is completely ruled out !
- ▶ to solve this: we assume that every word always appears at least once (or a smaller value)
 - ▶ ex: add-1 smoothing:

$$P(w_i | c_j) = \frac{\text{(frequency of } w_i \text{ in } c_j) + 1}{\text{total number of words in } c_j + \text{size of vocabulary}}$$

Add-one Smoothing

- ▶ Pretend we have seen every n-gram at least once
- ▶ Intuitively:
 - ▶ $\text{new_count(n-gram)} = \text{old_count(n-gram)} + 1$
- ▶ The idea is to give a little bit of the probability space to unseen events

Add-one: Example

unsmoothed bigram counts (frequencies):

	<i>2nd word</i>									
	<i>I</i>	want	to	eat	Chinese	food	lunch	...	Total	
<i>I</i>	8	1087	0	13	0	0	0	0	C(<i>I</i>)=3437	
want	3	0	786	0	6	8	6	...	C(want)=1215	
to	3	0	10	860	3	0	12	...	C(to)=3256	
eat	0	0	2	0	19	2	52	...	C(eat)=938	
Chinese	2	0	0	0	0	120	1	...	C(Chinese)=213	
food	19	0	17	0	0	0	0	...	C(food)=1506	
lunch	4	0	0	0	0	1	0	...	C(lunch)=459	
...								...		
								...		
									N=10,000	

- ▶ Assume a vocabulary of 1616 (different) words
 - ▶ $V = \{a, \text{aardvark}, \text{aardwolf}, \text{aback}, \dots, I, \dots, \text{want}, \dots, \text{to}, \dots, \text{eat}, \text{Chinese}, \dots, \text{food}, \dots, \text{lunch}, \dots, \text{zoophyte}, \text{zucchini}\}$
 - ▶ $|V| = 1616$ words
- ▶ And a total of $N = 10,000$ bigrams (~word instances) in the training corpus

Add-one: Example

unsmoothed bigram counts:

	<i>1st word</i>	<i>I</i>	want	to	eat	Chinese	food	lunch	...	<i>Total</i>
<i>1st word</i>	<i>I</i>	8	1087	0	13	0	0	0		$C(I)=3437$
	want	3	0	786	0	6	8	6		$C(want)=1215$
	to	3	0	10	860	3	0	12		$C(to)=3256$
	eat	0	0	2	0	19	2	52		$C(eat)=938$
	Chinese	2	0	0	0	0	120	1		$C(Chinese)=213$
	food	19	0	17	0	0	0	0		$C(food)=1506$
	lunch	4	0	0	0	0	1	0		$C(lunch)=459$
	...									
										N=10,000

unsmoothed bigram conditional probabilities:

	<i>I</i>	want	to	eat	Chinese	food	lunch	...	<i>Total</i>
<i>I</i>	0.002 (8/3437)	0.316	0	0.0037	0	0	0		
want	0.0025	0	0.647 (786/1215)	0	0.0049 (6/1215)	0.0066	0.0049		
to	0.0009	0	0.0030	0.264	0.0009	0	0.0037		
eat	0	0	0.002	0	0.020	0.002	0.0554		
Chinese	0.0094	0	0	0	0	0.56	0.0047		
food	0.0126	0	0.0113	0	0	0	0		
lunch	0.0087	0	0	0	0	0.0002	0		
...									

note :

$$P(I|I) = \frac{8}{10\,000}$$

$$P(I|I) = \frac{8}{3\,437}$$

Add-one, more formally

$$P_{Add1}(w_1 w_2 \dots w_n) = \frac{C(w_1 w_2 \dots w_n) + 1}{N + B}$$

如上例中 $P(\text{III}) = (8+1) / (3437+1616)$

3437 N: size of the corpus

i.e. number of n-gram tokens in training corpus

1616 B: number of "bins"

i.e. number of different n-gram types

i.e. number of cells in the matrix

e.g. for bigrams, it's (size of the vocabulary)²

Add-delta Smoothing

- ▶ every previously unseen n-gram is given a low probability
- ▶ but there are so many of them that too much probability mass is given to unseen events
- ▶ instead of adding 1, add some other (smaller) positive value δ

$$P_{Add}(w_1 w_2 \dots w_n) = \frac{C(w_1 w_2 \dots w_n) + \delta}{N + \delta B}$$

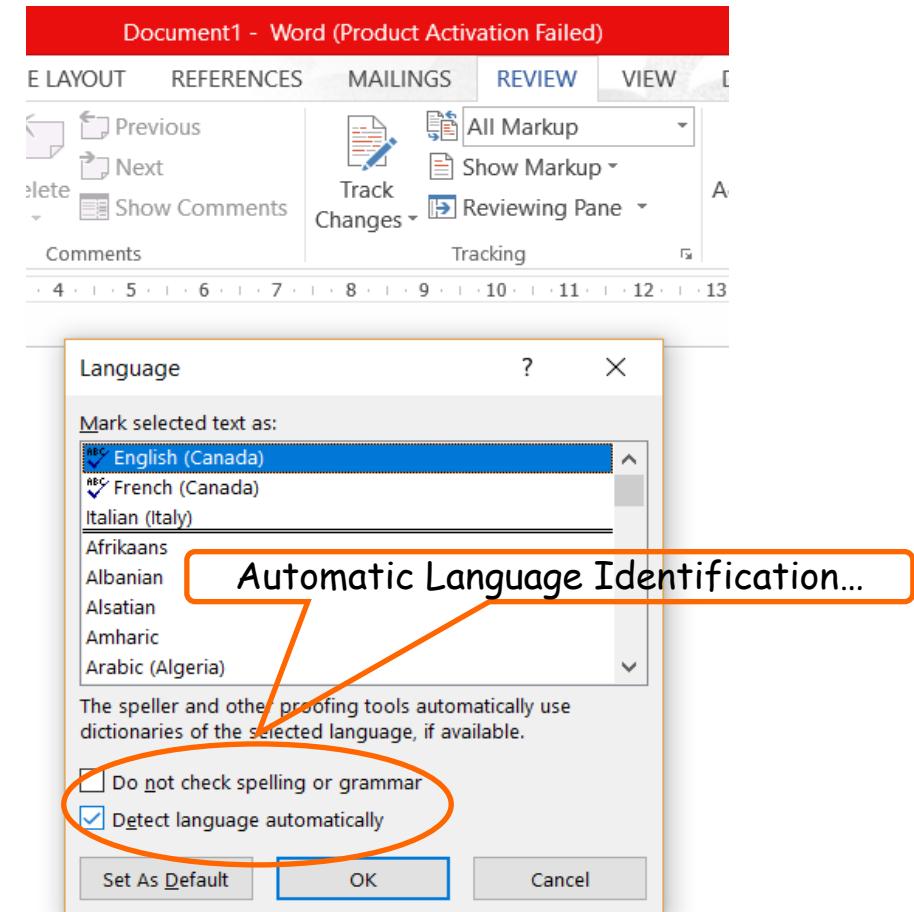
- ▶ most widely used value for $\delta = 0.5$
- ▶ better than add-one, but still...

Factors of Training Corpus

- ▶ Size:
 - ▶ the more, the better
 - ▶ but after a while, not much improvement...
 - ▶ bigrams (characters) after 100's million words
 - ▶ trigrams (characters) after some billions of words
- ▶ Genre (adaptation):
 - ▶ training on cooking recipes and testing on aircraft maintenance manuals

Example: Language Identification

- ▶ hypothesis: texts that resemble each other (same author, same language) share similar character/word sequences
 - ▶ In English character sequence “ing” is more probable than in French
- ▶ Training phase:
 - ▶ construction of the language model
 - ▶ with pre-classified documents (known language/author)
- ▶ Testing phase:
 - ▶ apply language model to unknown text



Example: Language Identification

- ▶ bigram of characters
 - ▶ characters = 26 letters (case insensitive)
 - ▶ possible variations: case sensitivity, punctuation, beginning/end of sentence marker, ...

Example: Language Identification

1. Train a character-based language model for Italian:

2. Train a character-based language model for Spanish:

Example: Language Identification

3. Given a unknown sentence “che bella cosa” is it in Italian or in Spanish?
 $P(\text{"che bella cosa"})$ with the Italian LM
 $P(\text{"che bella cosa"})$ with the Spanish LM
4. Highest probability -->language of sentence

Google's Web 1T 5-gram Model

- ▶ 5-grams
- ▶ generated from 1 trillion words
- ▶ 24 GB compressed
 - ▶ Number of tokens: 1,024,908,267,229
 - ▶ Number of sentences: 95,119,665,584
 - ▶ Number of unigrams: 13,588,391
 - ▶ Number of bigrams: 314,843,401
 - ▶ Number of trigrams: 977,069,902
 - ▶ Number of fourgrams: 1,313,818,354
 - ▶ Number of fivegrams: 1,176,470,663
- ▶ See discussion: <http://googleresearch.blogspot.com/2006/08/all-our-n-gram-are-belong-to-you.html>
- ▶ See Google N-gram Viewer: http://en.wikipedia.org/wiki/Google_Ngram_Visualizer

Problem with N-gram

- ▶ Natural language is not linear
- ▶ there may be *long-distance dependencies*.
 - ▶ Syntactic dependencies
 - ▶ The man next to the large oak tree near ... is tall.
 - ▶ The men next to the large oak tree near ... are tall.
 - ▶ Semantic dependencies
 - ▶ The bird next to the large oak tree near ... flies rapidly.
 - ▶ The man next to the large oak tree near ... talks rapidly.
 - ▶ World knowledge
 - ▶ Michael Jackson, who was featured in ..., is buried in California.
 - ▶ Michael Bublé, who was featured in, is living in California.
 - ▶ ...
- ▶ More complex models of language are needed to handle such dependencies.

Topics

- ▶ ***Word Embeddings***
- ▶ ***N-gram Models***
- ▶ Word Embeddings (continued)
- ▶ Contextualized Word Embeddings

Word2Vec

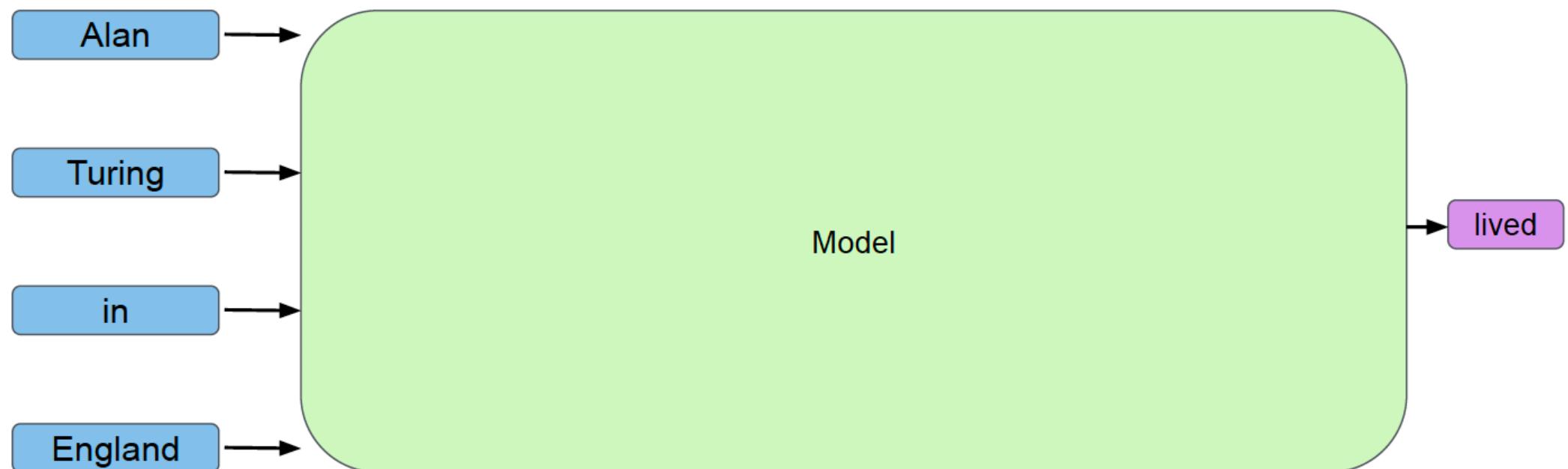
- ▶ Word2Vec is an efficient way to compute word embeddings.
- ▶ It is inspired from the distributional hypothesis.
- ▶ It is based on a very simple linear model.
- ▶ There are two versions:
 - ▶ Continuous bag-of-words.
 - ▶ Continuous Skip-Gram.

Word2Vec – Continuous bag-of-words

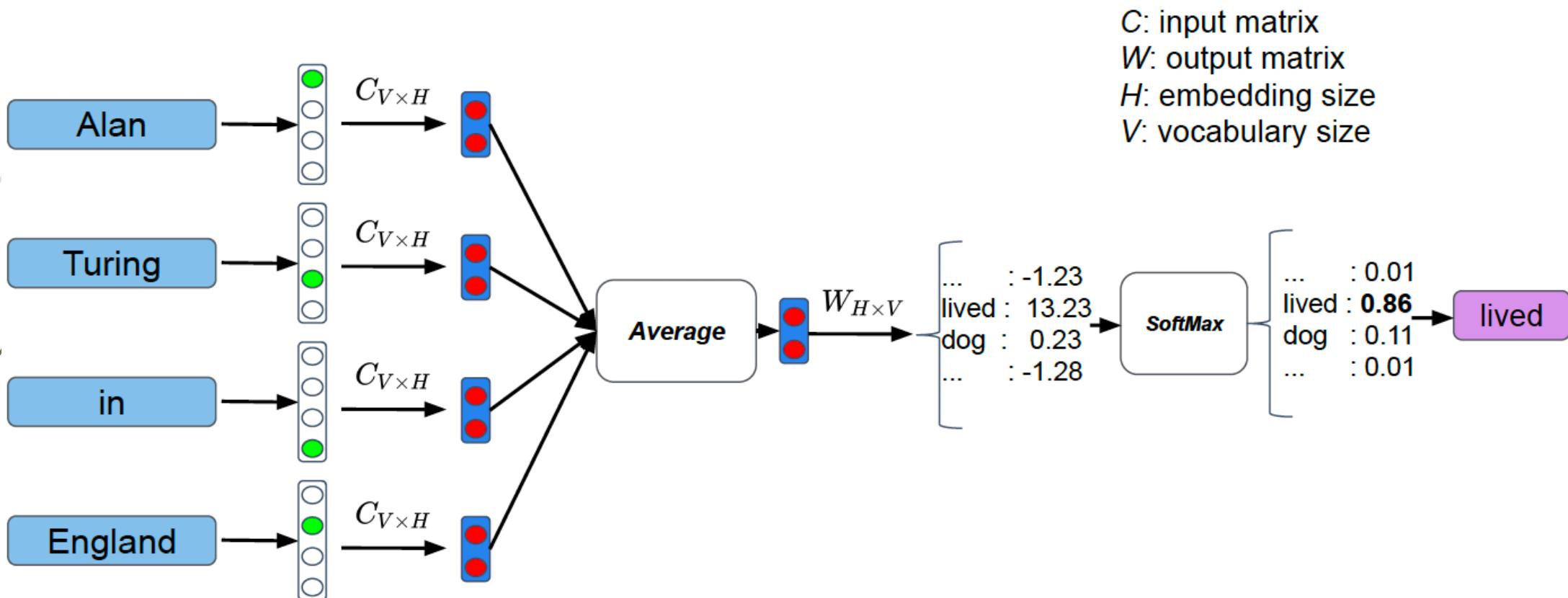
- ▶ The task is simple:
 - ▶ Mask a word in a sequence of text.
 - ▶ Train a model to predict this word given the words that surround it.
- ▶ For example, for “Alan Turing lived in England”:
 - ▶ Input: “Alan Turing in England”
 - ▶ Target: “lived”

Word2Vec – Continuous bag-of-words

- Task: predict a word given some context window.

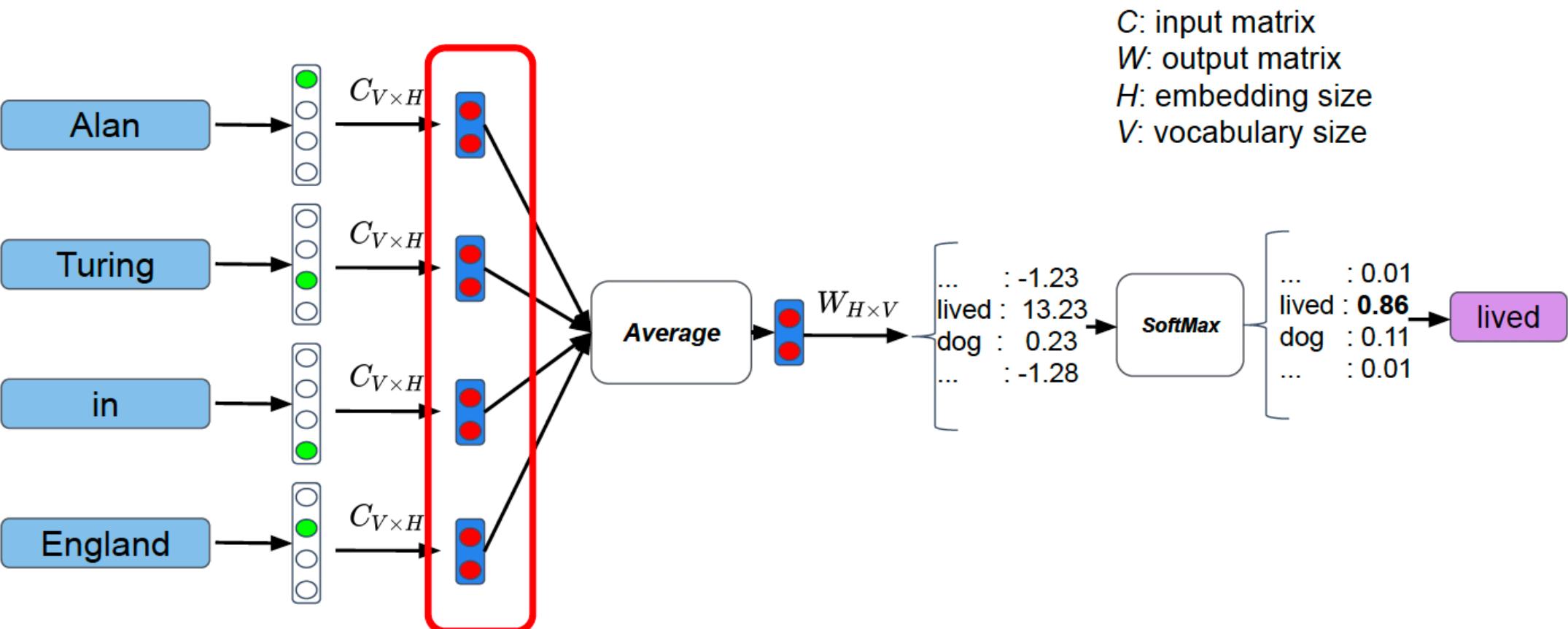


Word2Vec – Continuous bag-of-words



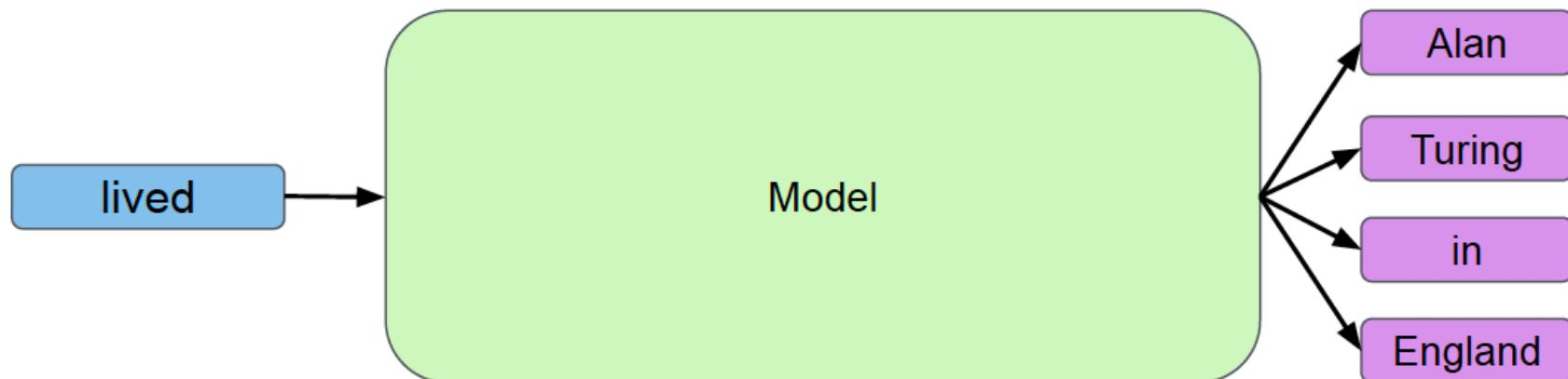
Word2Vec – Continuous bag-of-words

- The embeddings is what we really care about.



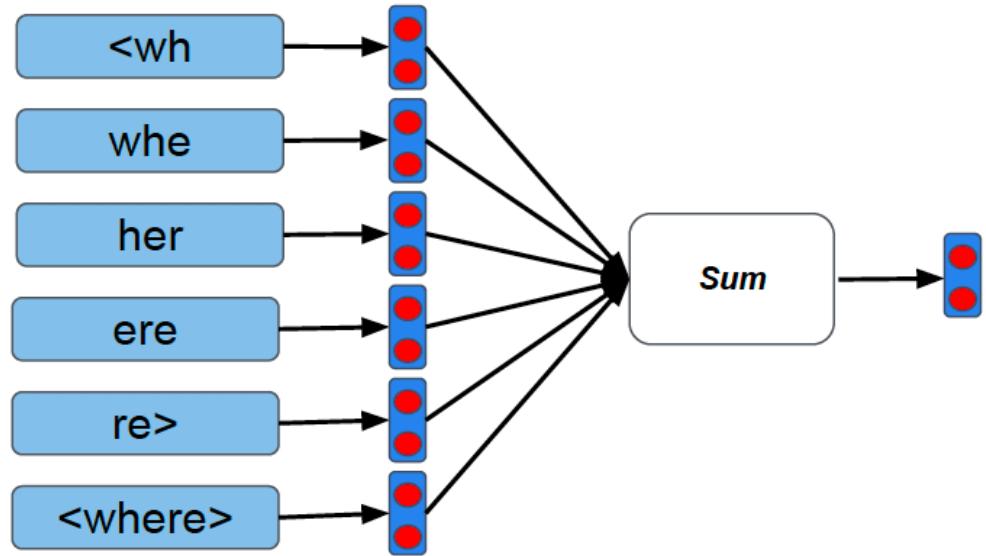
Word2Vec – Continuous Skip-Gram

- ▶ In the Continuous Skip-Gram version, the Word2Vec algorithm has a different goal:
 - ▶ It predicts the context given the central word.
 - ▶ It works better than continuous bag-of-words in the presence of smaller amounts of data.



FastText

- Extension of the Continuous Skip-Gram model to:
 - better handle out-of-vocabulary words,
 - consider the morphology of the words.
- The final embedding for one word is the sum of the original word embedding and the embeddings of its n-grams.



FastText representation of the word “where” if we only use 3-grams

Pre-Training

- ▶ Prior to training a model on a task of interest, the embeddings can be learned in a pre-training phase.
- ▶ Pre-training can be greatly beneficial given that algorithms like Word2Vec and FastText can learn embeddings in an unsupervised way, thus making it possible to leverage huge amounts of data.

Pre-Training

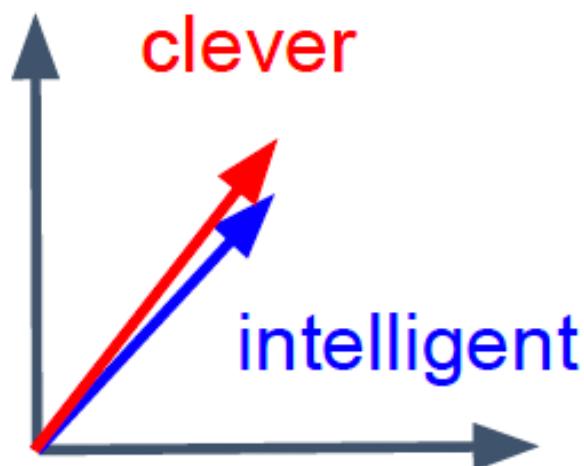
- ▶ Once pre-training is done, the embeddings can then be used for the real task of interest (i.e., the training phase).
 - ▶ This task is usually supervised and cannot therefore leverage large amounts of data.
 - ▶ This idea greatly helps to deal with scarcity of data in the training phase.

Topics

- ▶ ***Word Embeddings***
- ▶ ***N-gram Models***
- ▶ ***Word Embeddings (continued)***
- ▶ Contextualized Word Embeddings

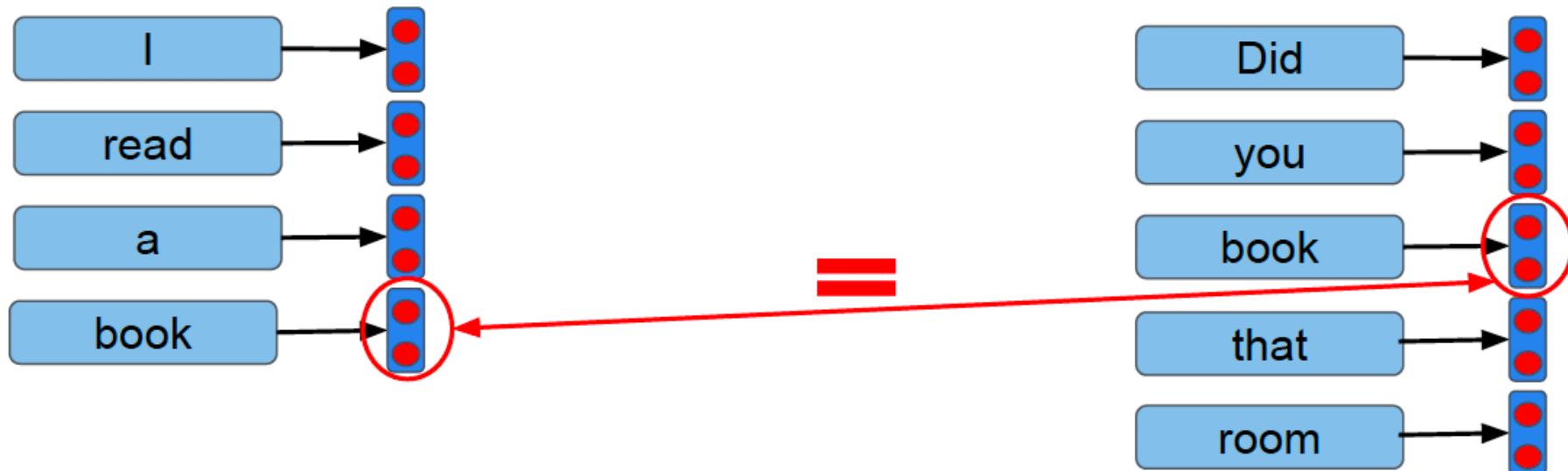
Word Embeddings and Synonymy

- Word embeddings are very useful representations that often lead to improved results.
- This is in part due to the fact that they can capture synonymy properties.
 - Synonyms in general have very close embedding vectors.



Word Embeddings and Polysemy

- Word embeddings do **not** help with polysemy though.
- A word such as “book” has the **same** embedding **regardless of the context**.



Word Embeddings and Polysemy

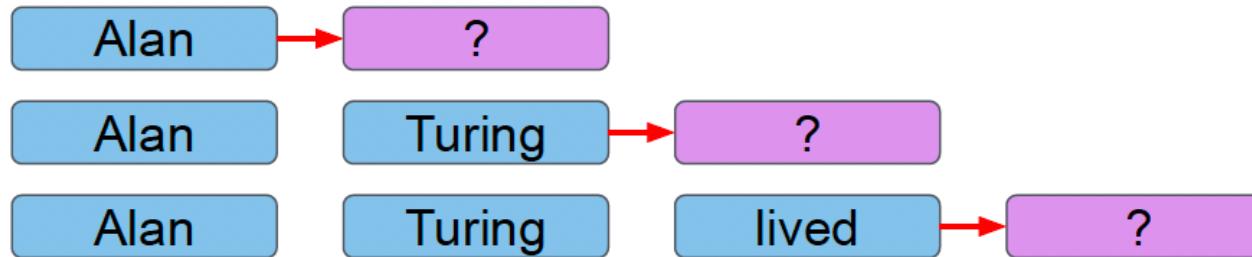
- ▶ Word embeddings reply to the question:
what is the embedding for “book”?
- ▶ In order to consider the context, the question should become:
what is the embedding for “book” in the sentence “I read a book”?

ELMo

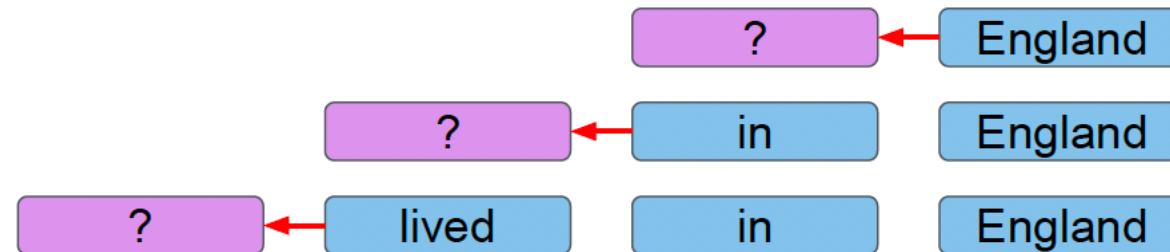
- ▶ ELMo (Embeddings from Language Models) proposes a solution to two problems:
 - ▶ Generate contextualized word embeddings...
 - ▶ ... in an unsupervised pre-training phase.
- ▶ To pre-train, they select the Language Modeling (LM) task.

Language Modeling

- Given n words (from a sentence), predict the next word ($n+1$).



- Language Modeling is an unsupervised task.
- Note it can also work 'right – to – left'.

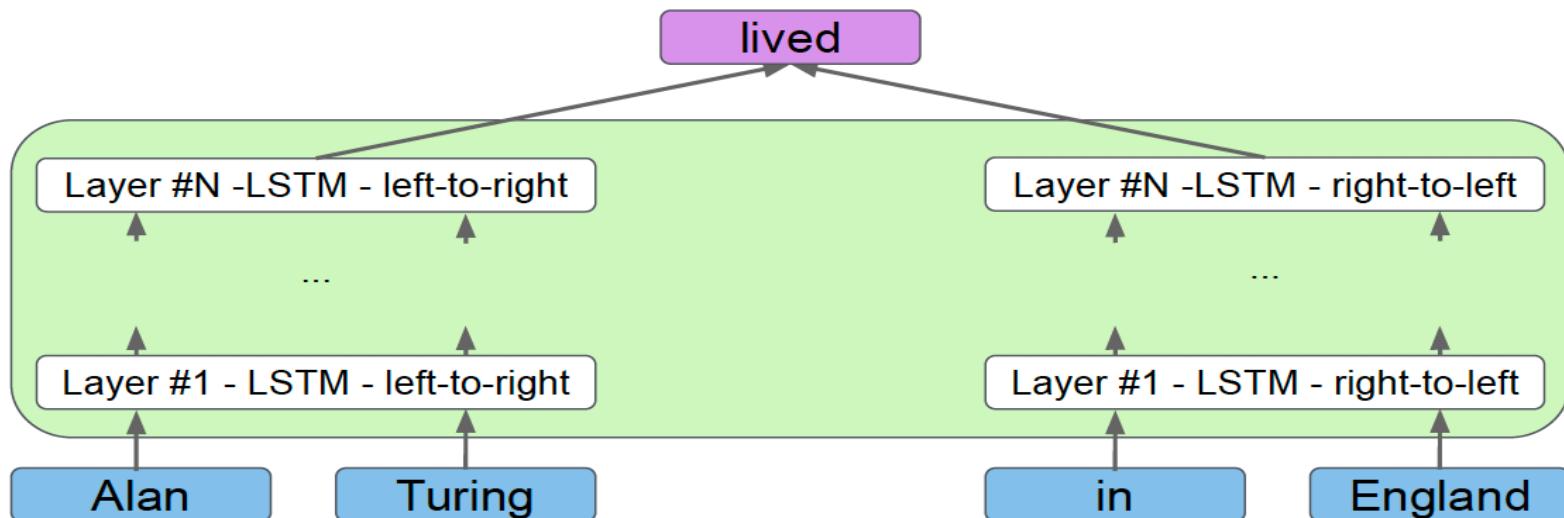


ELMo

- ▶ ELMo model is a N-layer bidirectional LSTM.

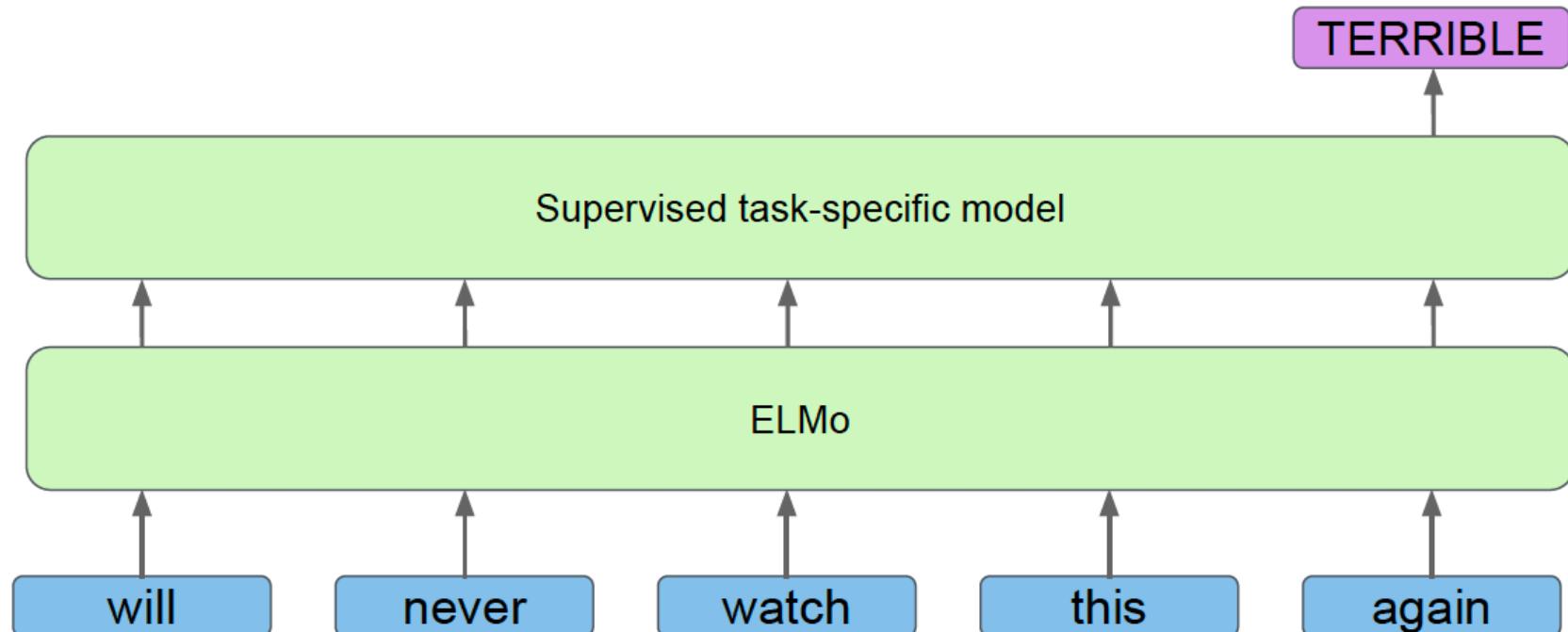
(Note: **LSTM** (Long Short Term Memory) Networks are called fancy recurrent neural networks with some additional features.)

- ▶ The contextualized embedding at a time step corresponds to a combination of the hidden states of all the various layers.

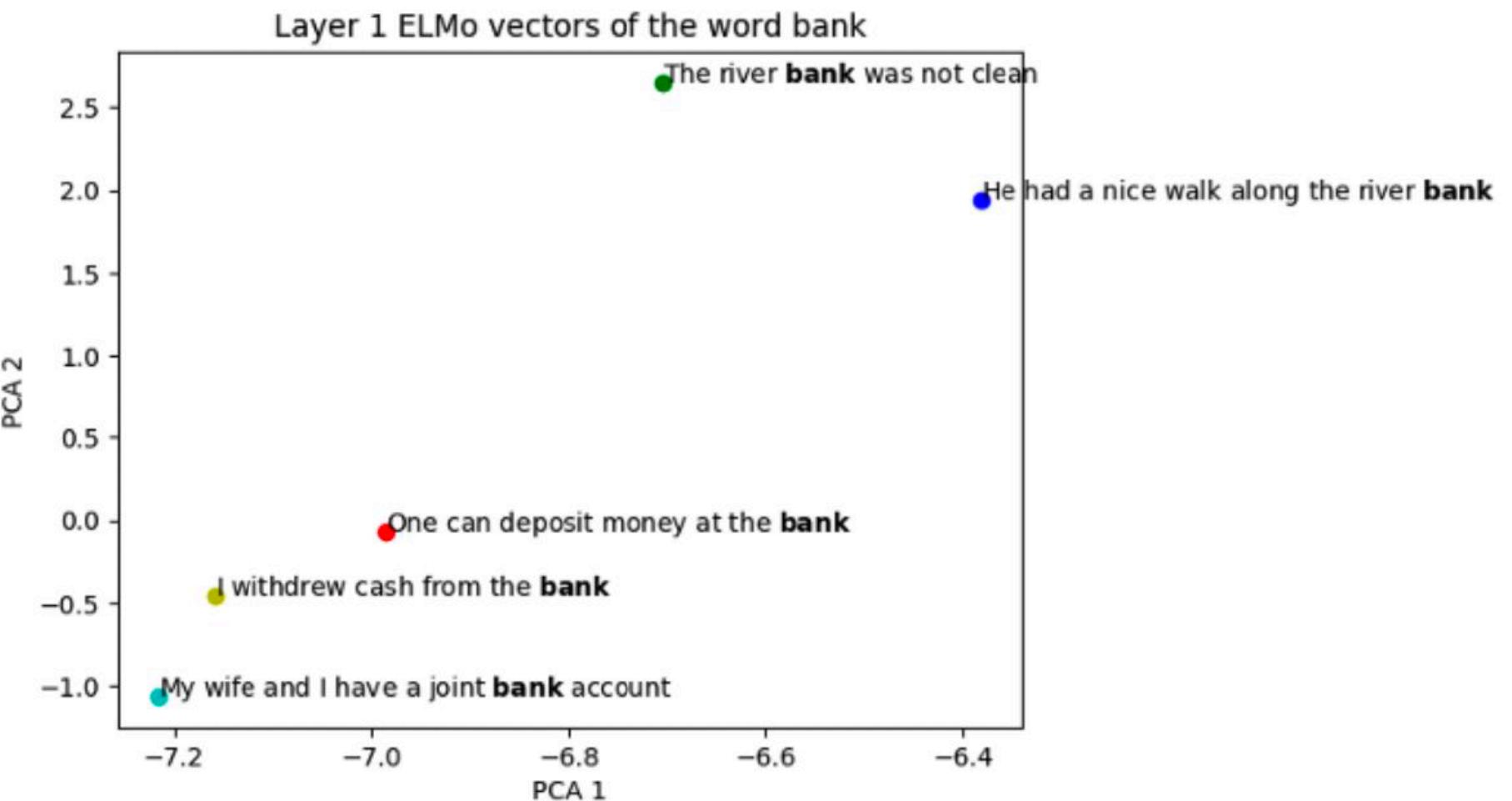


ELMo

- After pre-training, ELMo can be used with any other model (trained on supervised task).



ELMo - Visualization



The End

