

# Statistical Machine Translation

# What is machine translation?

نائب امریکی صدر ڈک چینی کا کہنا ہے کہ میں اسامہ بن لادن  
کو زندہ یا مردہ دیکھنا چاہتا ہوں۔

American Vice President Dick Cheney has said that he wants to see  
Osama bin Laden dead or alive.

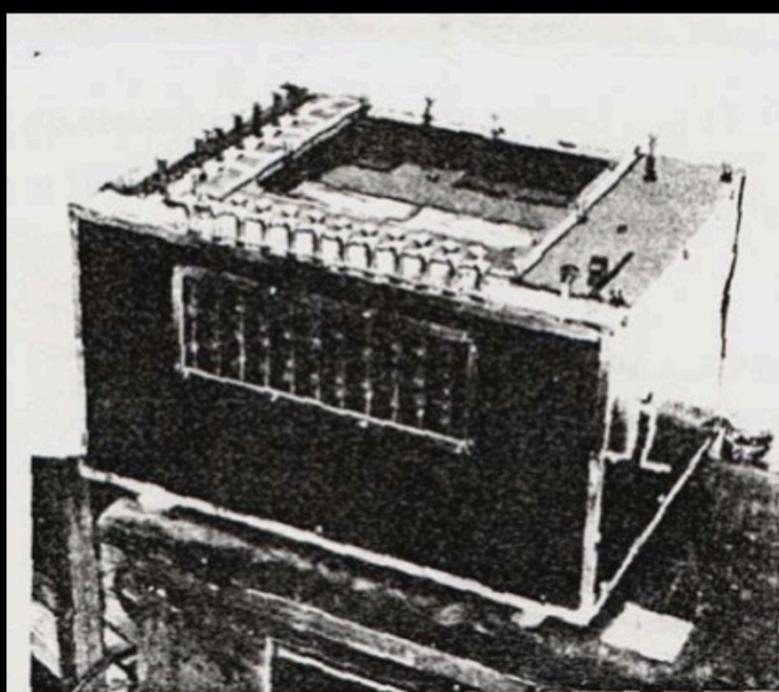
# A brief history



The Tower of Babel

Pieter Brueghel the Elder (1563)

# A brief history



mechanical brain  
(Artsrouni, France 1933)

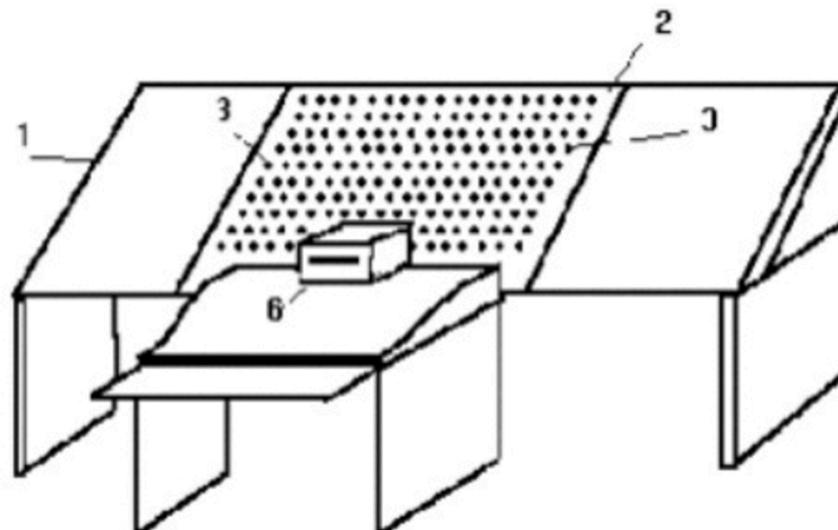
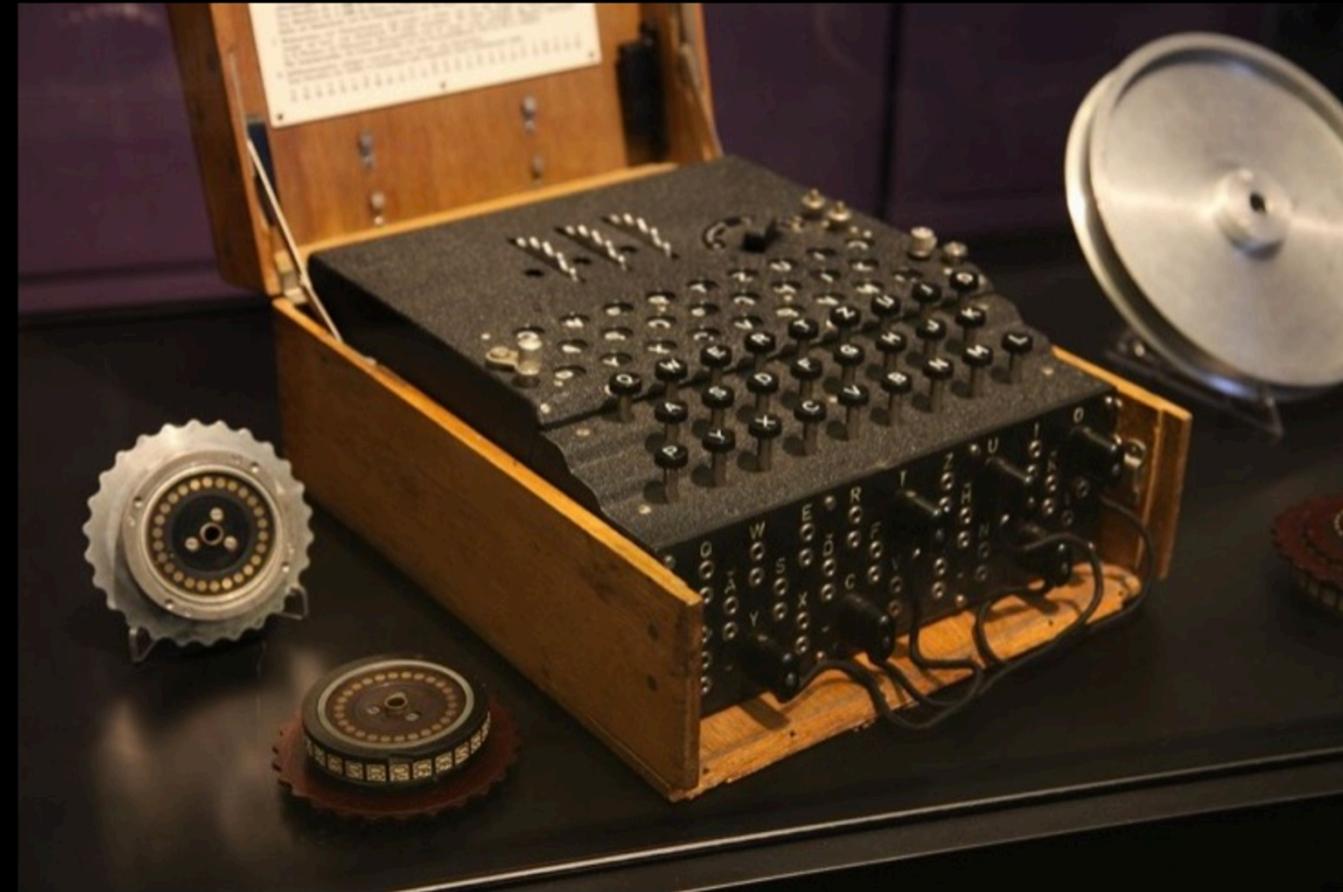


Fig. 3. Trojanskij's translating machine (*from Bel'skaja et al. 1959*)

(Trojanskij, Russia 1933)

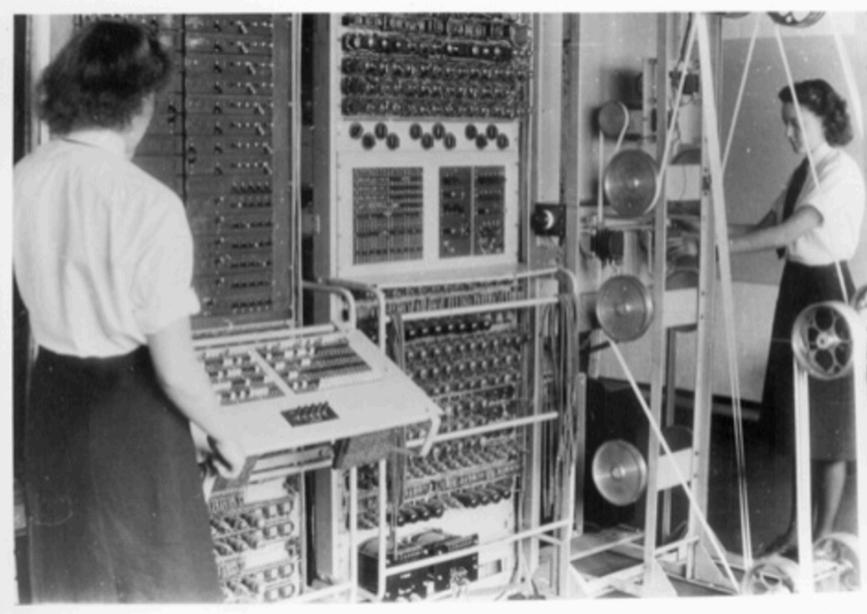
# A brief history

1940s

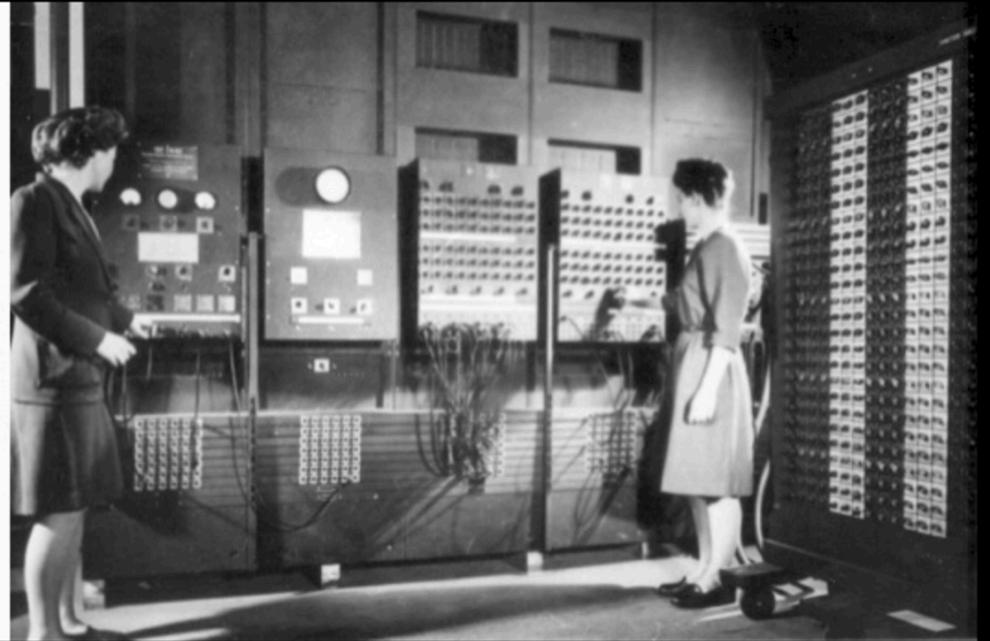


First NLP problem: the German Enigma

# A brief history



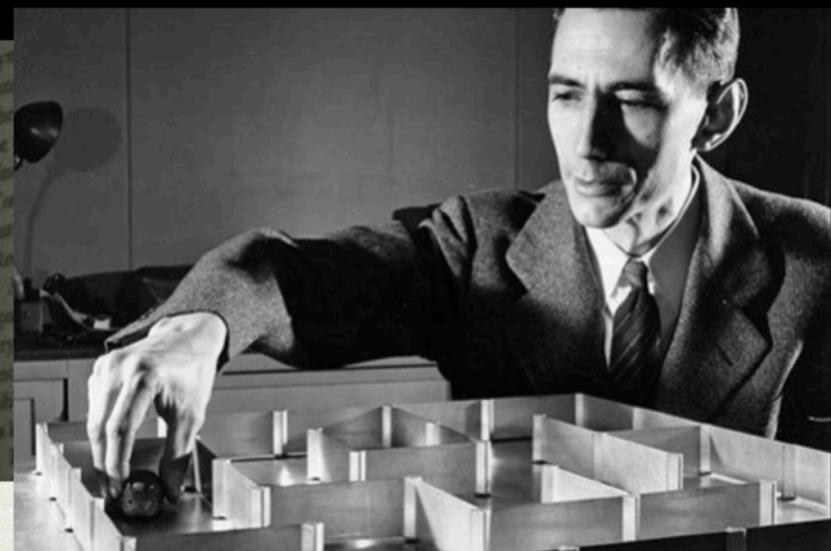
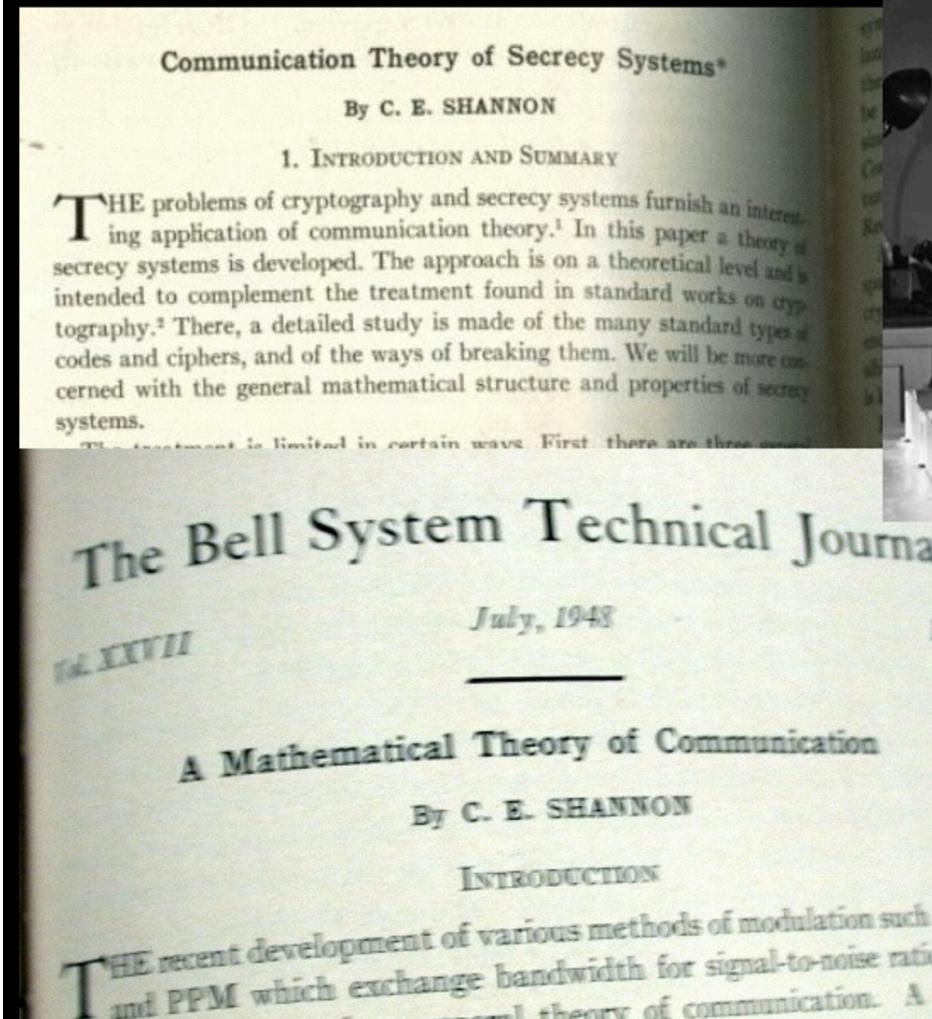
Colussus (1944)



ENIAC (1946)

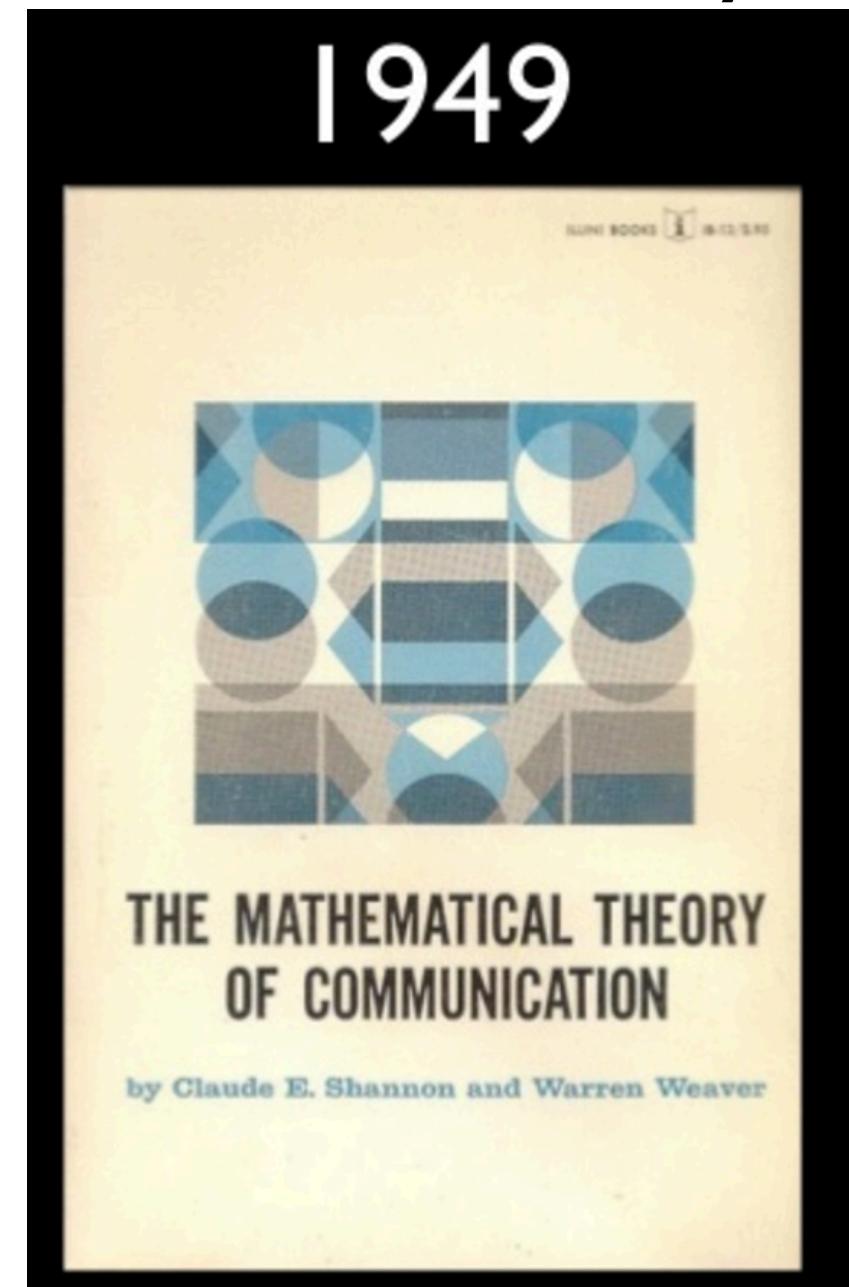
# A brief history

1948



# A brief history

1949



# A brief history



*When I look at an article in Russian, I say: "This is really written in English, but it has been coded in some strange symbols. I will now proceed to decode."*

Warren Weaver (1949)

# A brief history

1962



*Association for Machine Translation  
and Computational Linguistics  
founded*

# A brief history

1966

“‘Machine Translation’ presumably means going by algorithm from machine-readable source text to useful target text, without recourse to human translation or editing. In this context, there has been no machine translation of general scientific text, and none is in immediate prospect.”

*The ALPAC report*

# A brief history

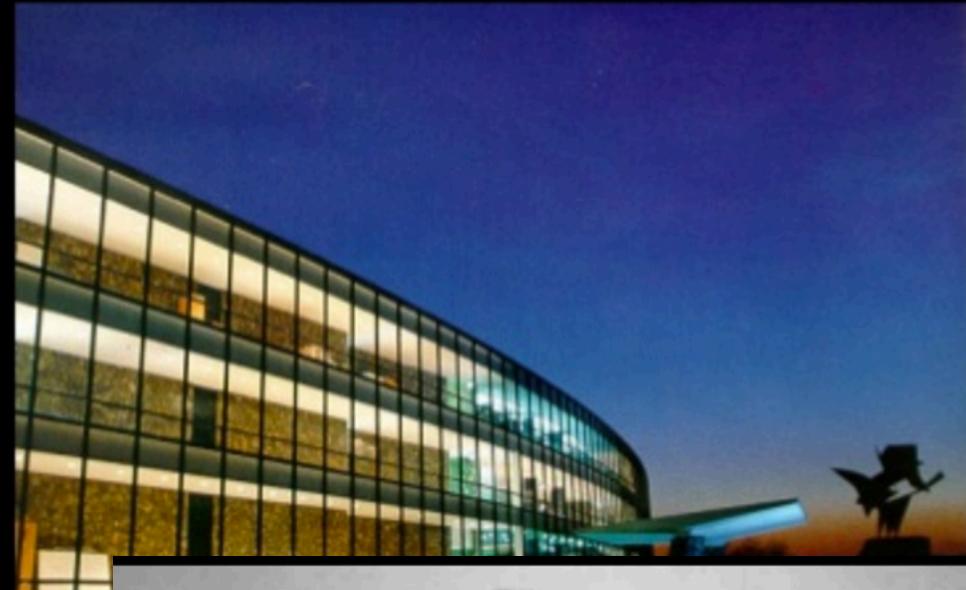
1968



*Association for Machine Translation and  
Computational Linguistics*  
becomes  
*Association for Computational Linguistics*

# A brief history

1972



# A brief history

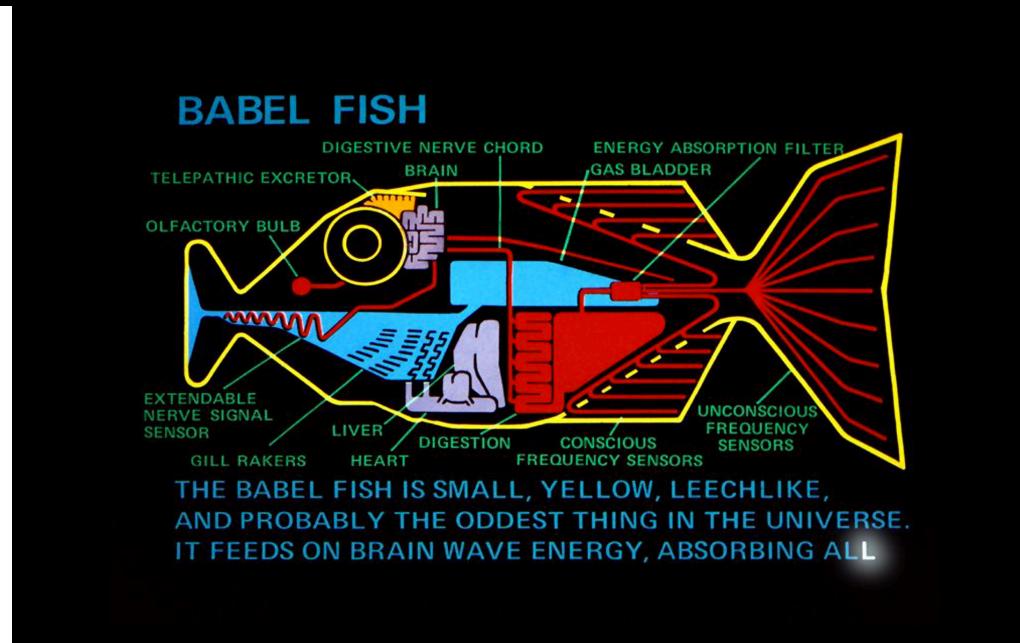
1975

$$\hat{W} = \arg \max_W P(W|A) = \arg \max_W P(A|W)P(W)$$

“Here’s something we jocularly called the fundamental equation of speech recognition at ICASSP in 1981. As I’m sure all of you very well know, this is just an application of Bayes’ Rule.”

*-Bob Mercer*

# A brief history



# A brief history

1990

## A STATISTICAL APPROACH TO MACHINE TRANSLATION

Peter F. Brown, John Cocke, Stephen A. Della Pietra, Vincent J. Della Pietra, Fredrick Jelinek,  
John D. Lafferty, Robert L. Mercer, and Paul S. Roossin

IBM

Thomas J. Watson Research Center  
Yorktown Heights, NY

In this paper, we present a statistical approach to machine translation. We describe the application of our approach to translation from French to English and give preliminary results.

### 1 INTRODUCTION

The field of machine translation is almost as old as the modern digital computer. In 1949 Warren Weaver suggested that the problem be attacked with statistical methods and ideas from information theory, an area which he, Claude Shannon, and others were developing at the time (Weaver 1949). Although researchers quickly abandoned this approach, advancing numerous theoretical objections, we believe that the true obstacles lay in the relative impotence of the available computers and the dearth of machine-readable text from which to gather the statistics vital to such an attack. Today, computers are five orders of magnitude faster than they were in 1950 and have hundreds of millions of bytes of storage. Large, machine-readable corpora are readily available. Statistical methods have proven their value in automatic speech recognition (Bahl et al. 1983) and have recently been applied to lexicography (Sinclair 1985) and to natural language processing (Baker

sentence in one language is a possible translation of any sentence in the other. We assign to every pair of sentences ( $S, T$ ) a probability,  $\Pr(T|S)$ , to be interpreted as the probability that a translator will produce  $T$  in the target language when presented with  $S$  in the source language. We expect  $\Pr(T|S)$  to be very small for pairs like (*Le matin je me brosse les dents* | *President Lincoln was a good lawyer*) and relatively large for pairs like (*Le président Lincoln était un bon avocat* | *President Lincoln was a good lawyer*). We view the problem of machine translation then as follows. Given a sentence  $T$  in the target language, we seek the sentence  $S$  from which the translator produced  $T$ . We know that our chance of error is minimized by choosing that sentence  $S$  that is most probable given  $T$ . Thus, we wish to choose  $S$  so as to maximize  $\Pr(S|T)$ . Using Bayes' theorem, we can write

$$\Pr(S|T) = \frac{\Pr(S) \Pr(T|S)}{\Pr(T)}$$

# A brief history



## Statistical Machine Translation Live

4/28/2006  
Franz Och

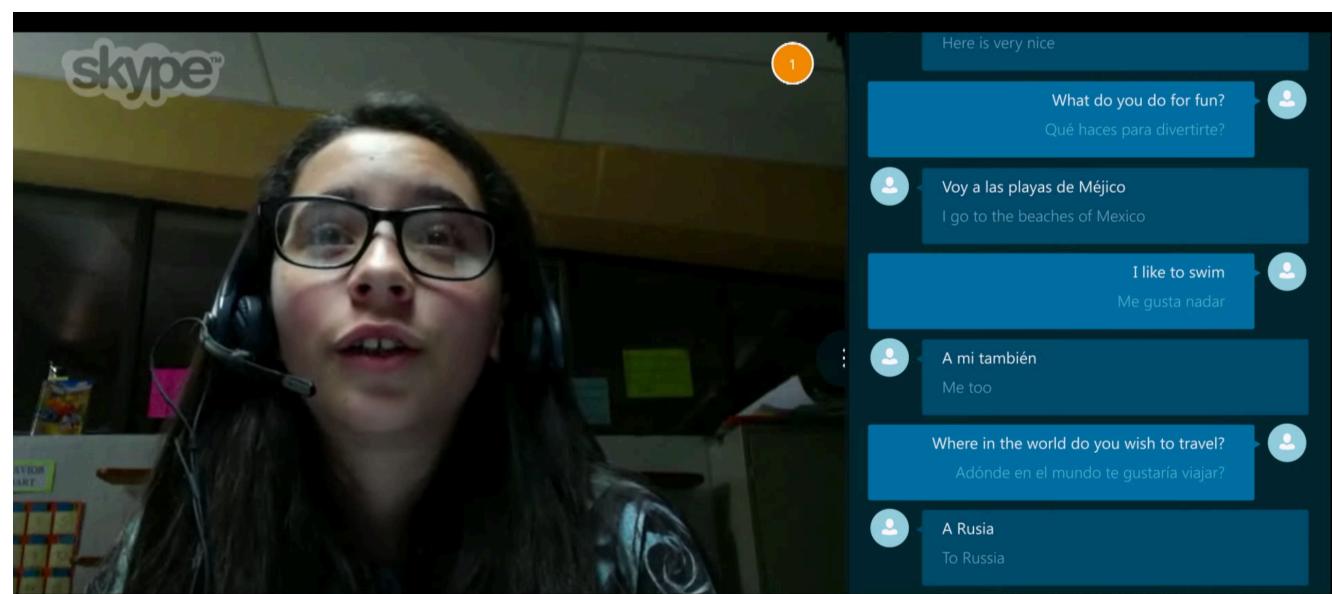
Because we want to provide everyone with access to all the world's information, including information written in every language, one of the exciting projects at Google Research is machine translation... Now you can see the results for yourself. We recently launched an online version of our system for Arabic-English and English-Arabic. Try it out!

# A brief history

Google Translate

English	Spanish	French	Detect language	▼	↔	English	Spanish	Arabic	▼
Afrikaans	Cebuano	Finnish	Hungarian	Latin	Romanian	Turkish			
Albanian	Chinese (Simplified)	French	Icelandic	Latvian	Russian	Ukrainian			
Arabic	Chinese (Traditional)	Galician	Indonesian	Lithuanian	Serbian	Urdu			
Armenian	Croatian	Georgian	Irish	Macedonian	Slovak	Vietnamese			
Azerbaijani	Czech	German	Italian	Malay	Slovenian	Welsh			
Basque	Danish	Greek	Japanese	Maltese	Spanish	Yiddish			
Belarusian	Dutch	Gujarati	Javanese	Marathi	Swahili				
Bengali	English	Haitian Creole	Kannada	Norwegian	Swedish				
Bosnian	Esperanto	Hebrew	Khmer	Persian	Tamil				
Bulgarian	Estonian	Hindi	Korean	Polish	Telugu				
Catalan	Filipino	Hmong	Lao	Portuguese	Thai				

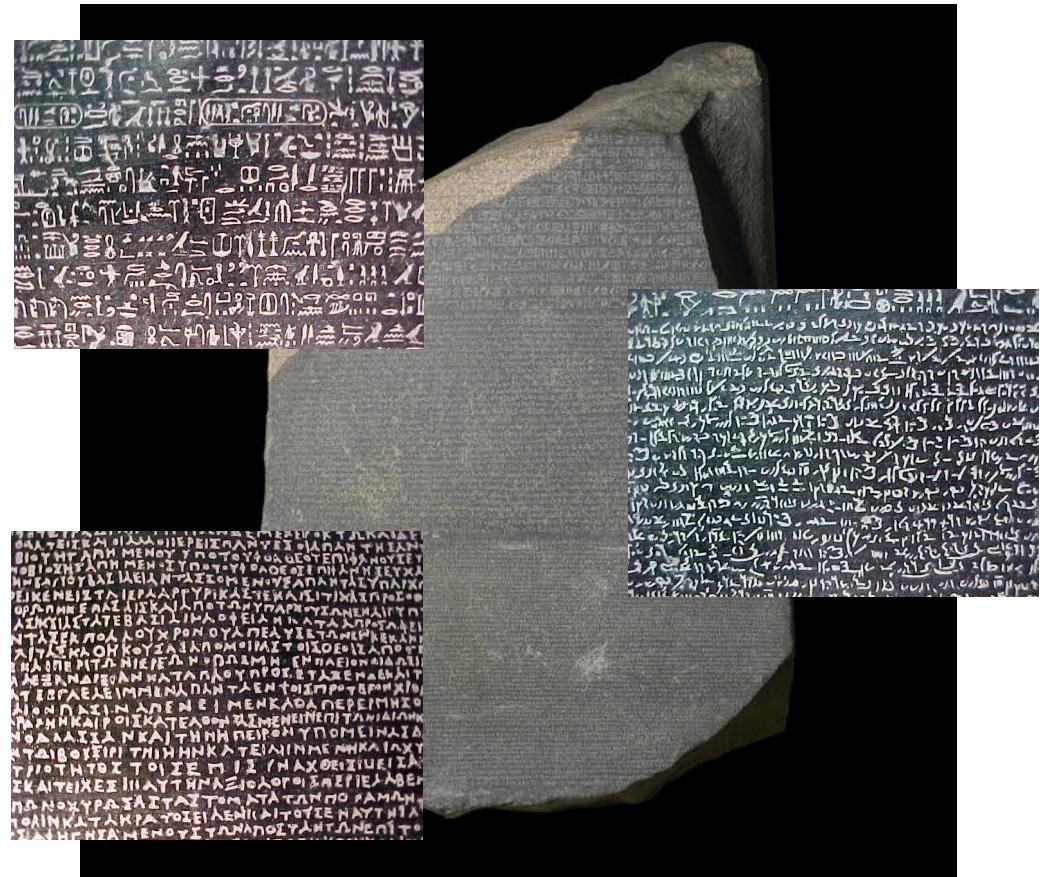
# A brief history



# The main ingredient



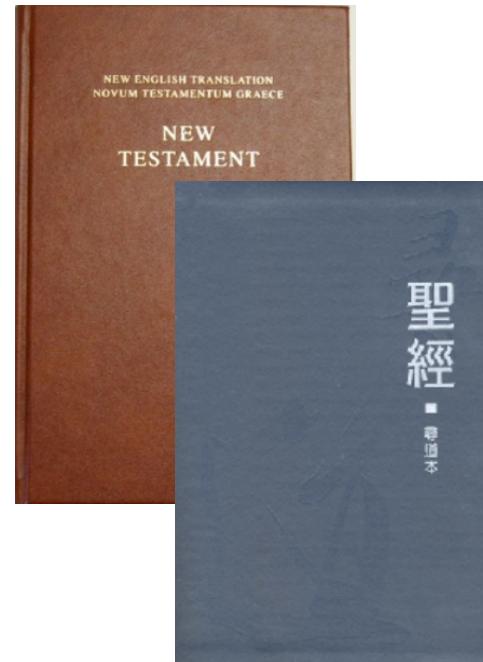
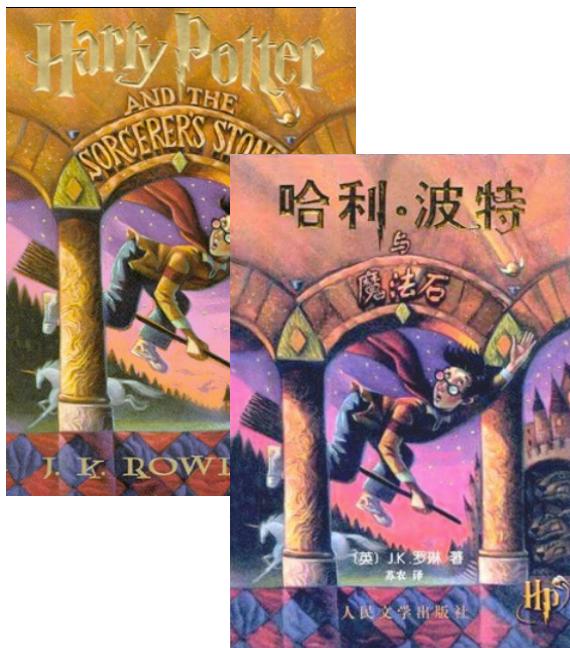
Jean-François Champollion  
Rosetta Stone translation, 1818



# The main ingredient



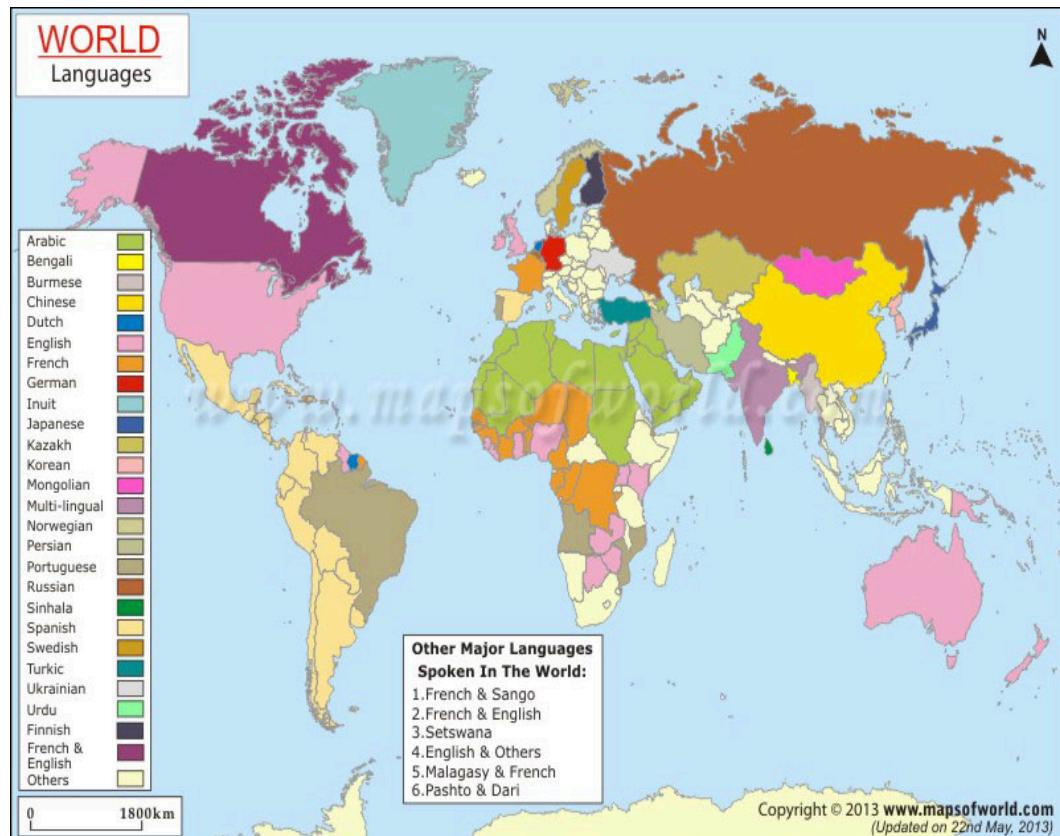
A screenshot of the UN website. The top navigation bar includes links for 'Peace and Security', 'Development', 'Human Rights', 'Humanitarian Affairs', and 'International Law'. The main content area features a large image of a woman in a blue headscarf and a headline about a Ban condemnation. On the right side, there are sections for 'Main Bodies' (General Assembly, Security Council, Economic & Social Council, etc.) and 'The UN and ...' (Sustainable food systems, World Food Day, etc.). The bottom right corner shows a banner for 'UNAMID' in Arabic.



# Interest in MT

## Languages in the world

- 6,800 living languages
- 600 with written tradition
- 100 languages are spoken by 95% of world population



## Translation market

- \$26 Billion global market (2010)
- Doubling every five years
- EU spends \$1 billion/year on Translation

## Academic

- A challenging task
- Boost other NLP tasks

# MT is hard

## A human translator needs ...

- ▶ to understand the source language
- ▶ to know how to speak the target language (well!)
- ▶ knowledge about the relationship between source and target language
- ▶ knowledge about the topic of the text to be translated
- ▶ knowledge about culture, values, traditions and expectation av speakers of source and target language

Computers have big problems with all of these issues

# MT is hard

## Word Translation Problems

- Words are ambiguous

He deposited money in a bank account  
with a high interest rate.

Sitting on the bank of the Mississippi,  
a passing ship piqued his interest.

- How do we find the right meaning, and thus translation?
- Context should be helpful

# MT is hard

## Syntactic Translation Problems

- Languages have different sentence structure

das	behaupten	sie	wenigstens
this	claim	they	at least
the		she	

- Convert from object-verb-subject (OVS) to subject-verb-object (SVO)
- Ambiguities can be resolved through syntactic analysis
  - the meaning **the** of **das** not possible (not a noun phrase)
  - the meaning **she** of **sie** not possible (subject-verb agreement)

# MT is hard

## Semantic Translation Problems

- Pronominal anaphora

I saw the movie and **it** is good.

- How to translate **it** into German (or French)?
  - **it** refers to **movie**
  - **movie** translates to **Film**
  - **Film** has masculine gender
  - ergo: **it** must be translated into masculine pronoun **er**

# MT is hard

## Semantic Translation Problems

- Coreference

Whenever I visit my uncle and his daughters,  
I can't decide who is my favorite cousin.

- How to translate **cousin** into German? Male or female?
- Complex inference required

# MT is hard

## Semantic Translation Problems

- Discourse

Since you brought it up, I do not agree with you.

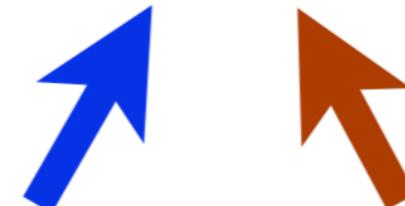
Since you brought it up, we have been working on it.

- How to translated **since**? Temporal or conditional?
- Analysis of discourse structure — a hard problem

# Two views of MT

- **Direct modeling** (aka pattern matching)
  - I have **really good learning algorithms** and a bunch of **example inputs** (source language sentences) and **outputs** (target language translations)
- **Code breaking** (aka the noisy channel, Bayes rule)
  - I know the **target language**
  - I have example **translations texts** (example enciphered data)

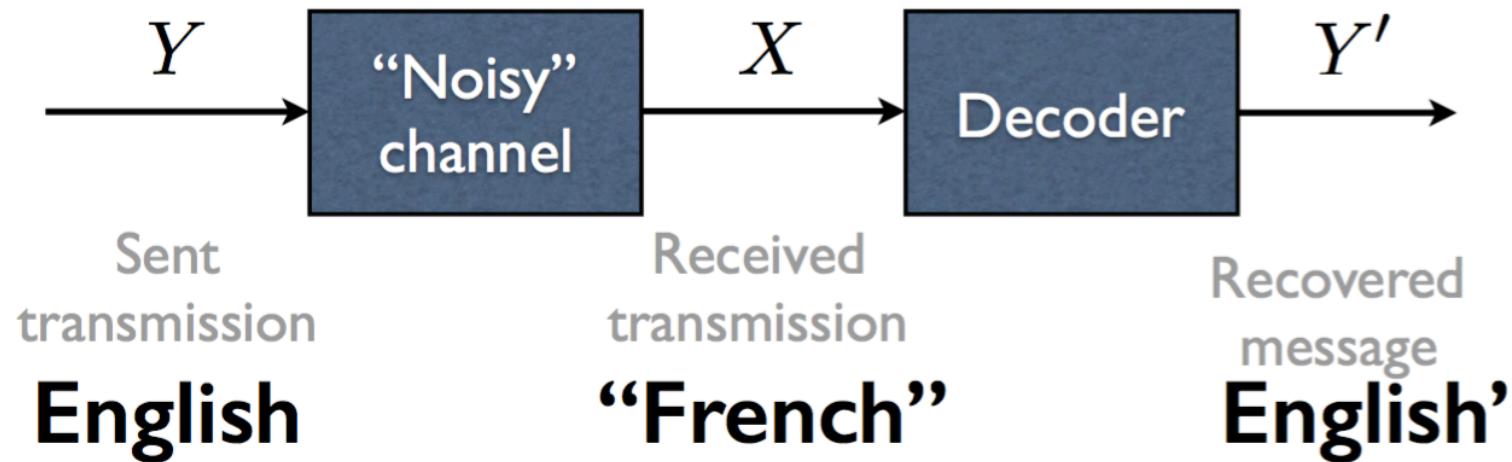
# MT as direct modeling

$$\hat{e} = \arg \max_e p_{\lambda}(e | f)$$


target      source

- one model does everything
- trained to reproduce a corpus of translations

# MT as code breaking (noisy channel)



$$\hat{e} = \arg \max_e p_{\varphi}(e) \times p_{\theta}(f | e)$$

language model

translation model

# Which view in 2016?

- Direct modeling is where most of the action is
  - Neural networks are very good at generalizing - and conceptually very simple
  - Inference in “product of two models” is hard
- Noisy channel ideas are incredibly important and still play a big role in how we think about translation
- We will cover both views

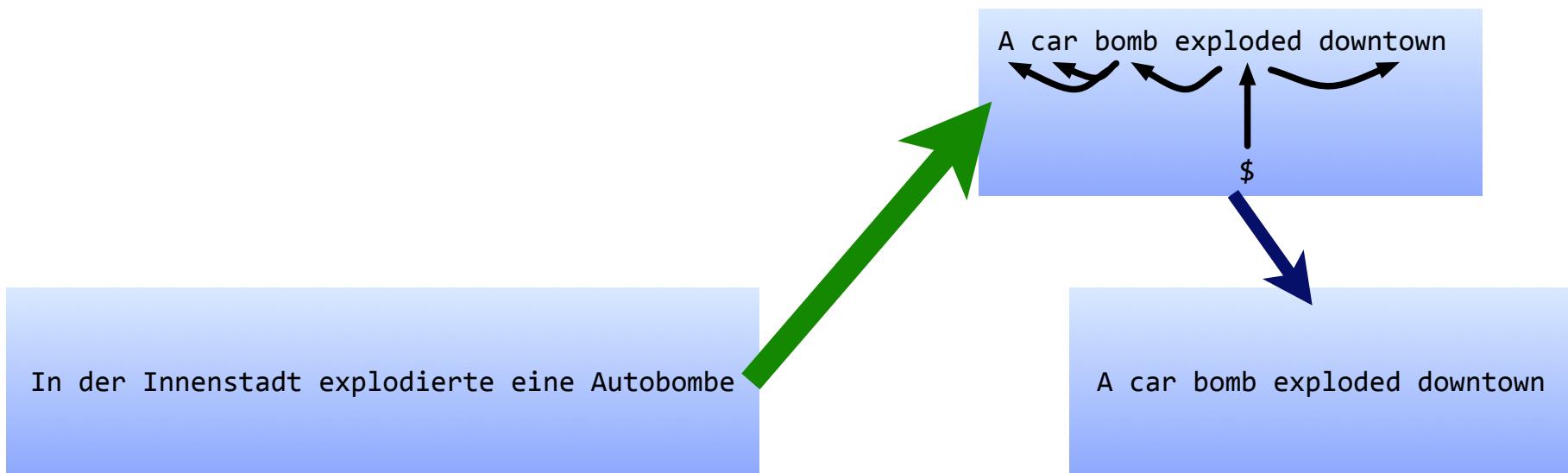
# Translation Modeling

In der Innenstadt explodierte eine Autobombe

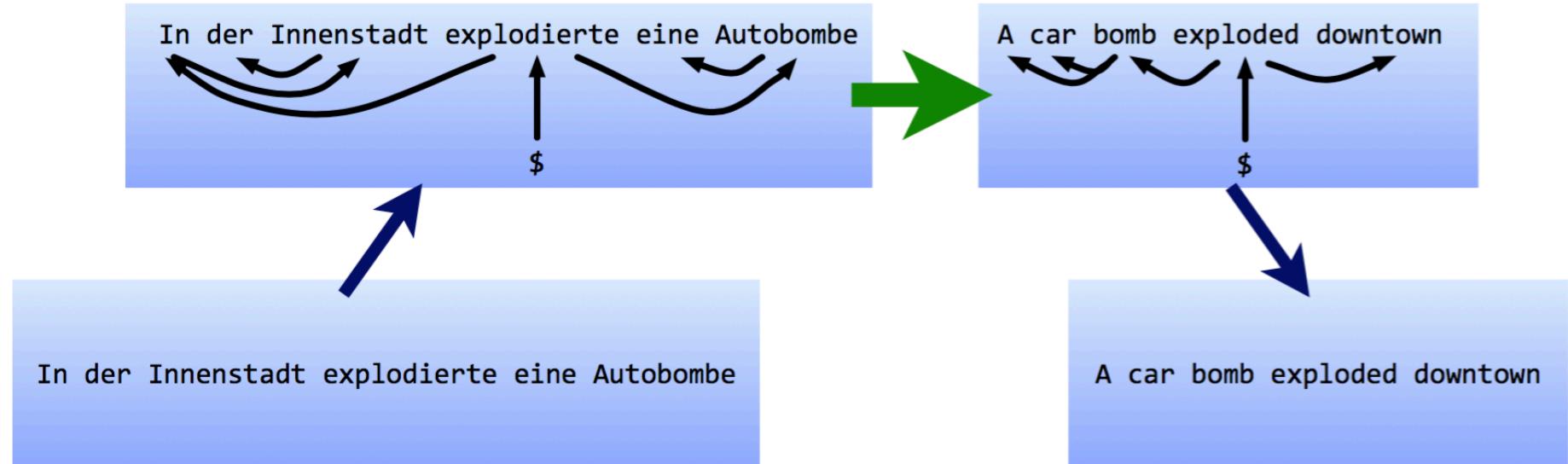


A car bomb exploded downtown

# Translation Modeling



# Translation Modeling



# Translation Modeling

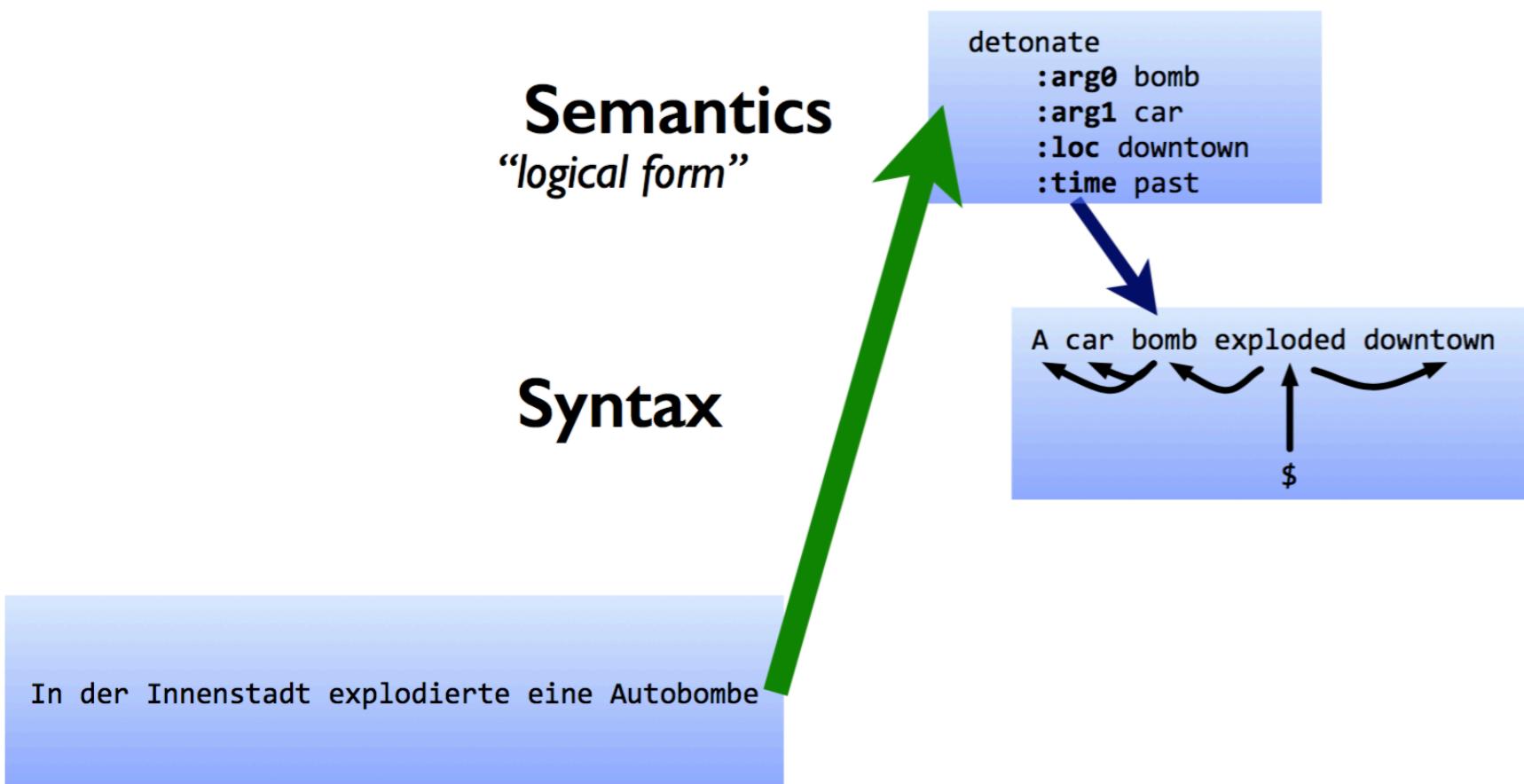
**Semantics**  
“logical form”

In der Innenstadt explodierte eine Autobombe

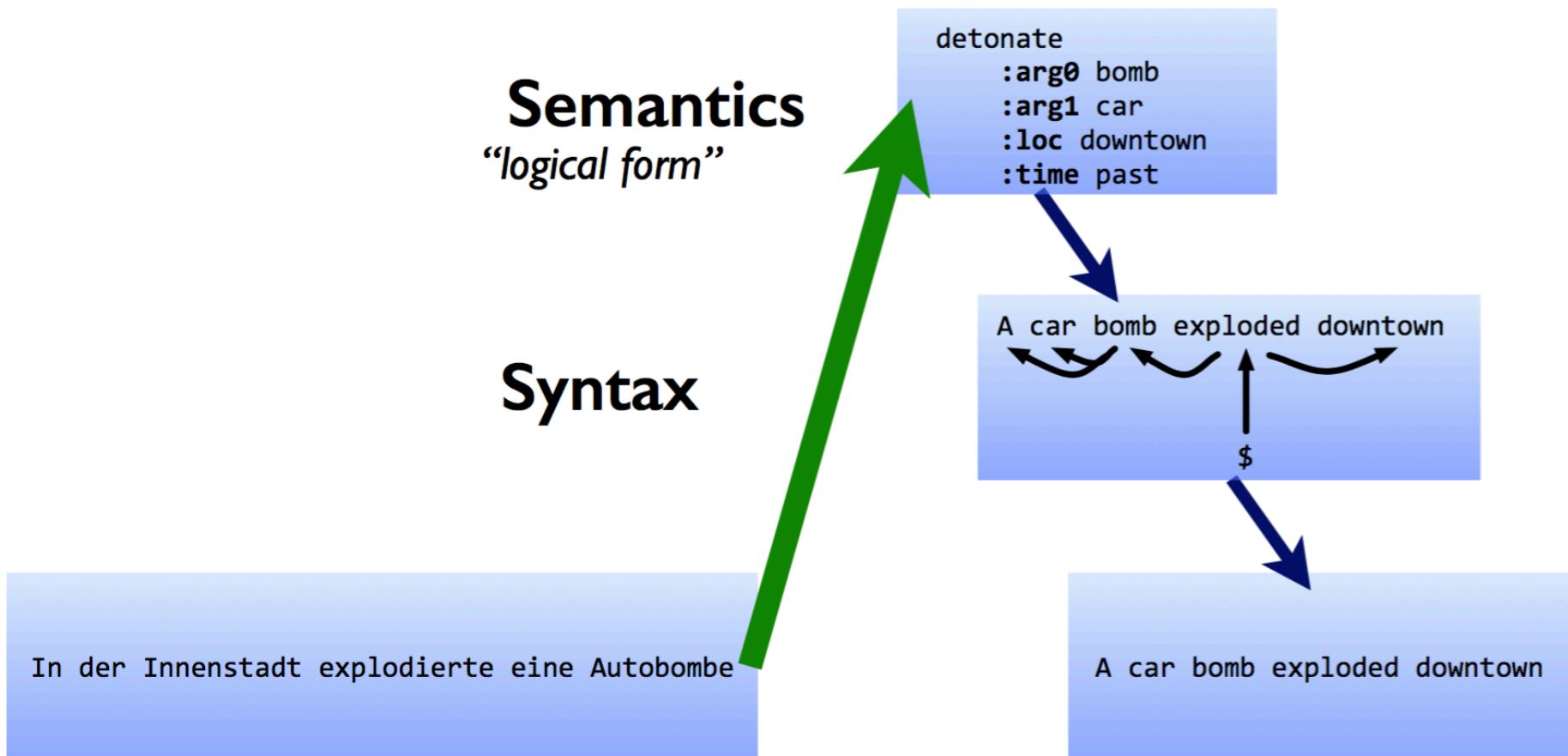
detonate  
:arg0 bomb  
:arg1 car  
:loc downtown  
:time past



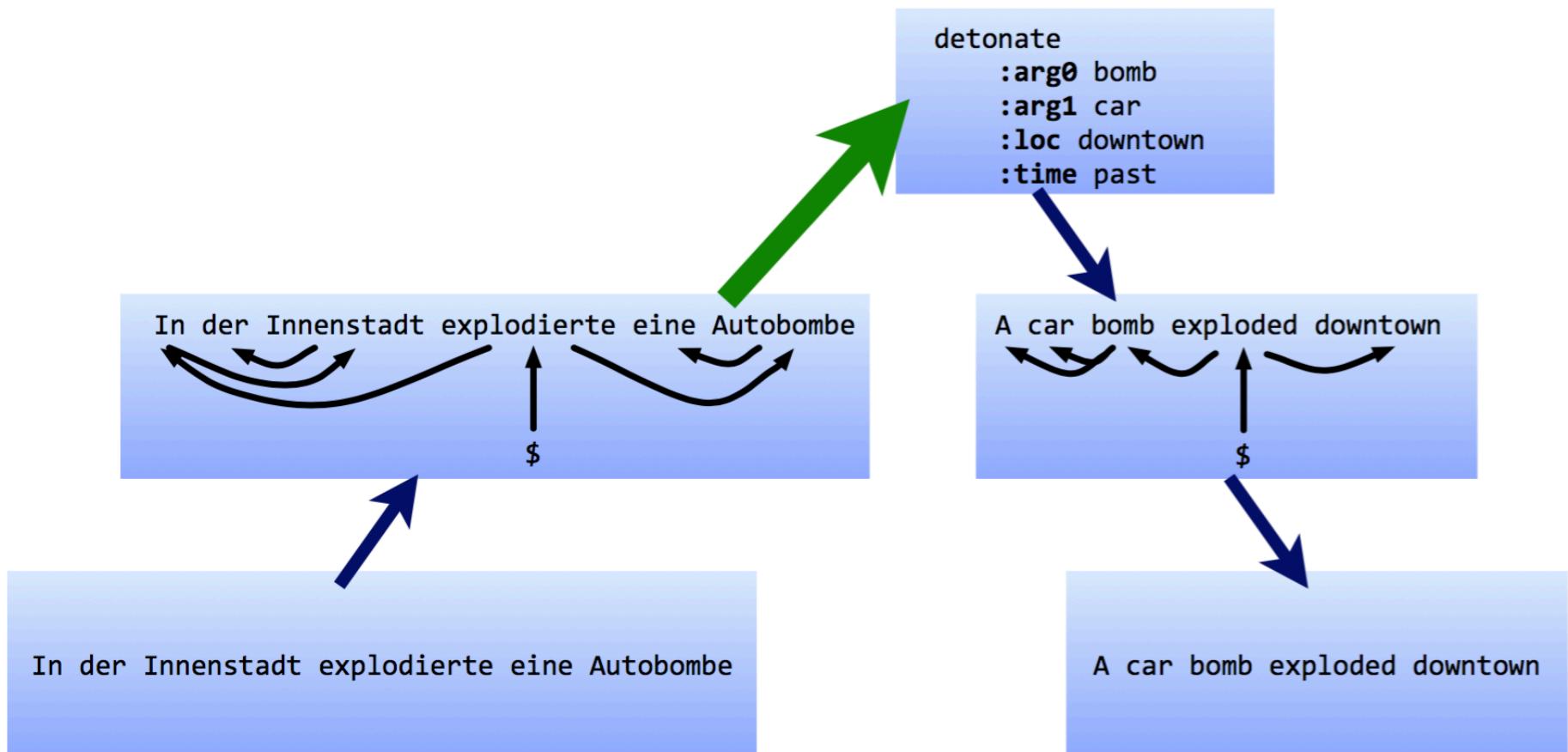
# Translation Modeling



# Translation Modeling

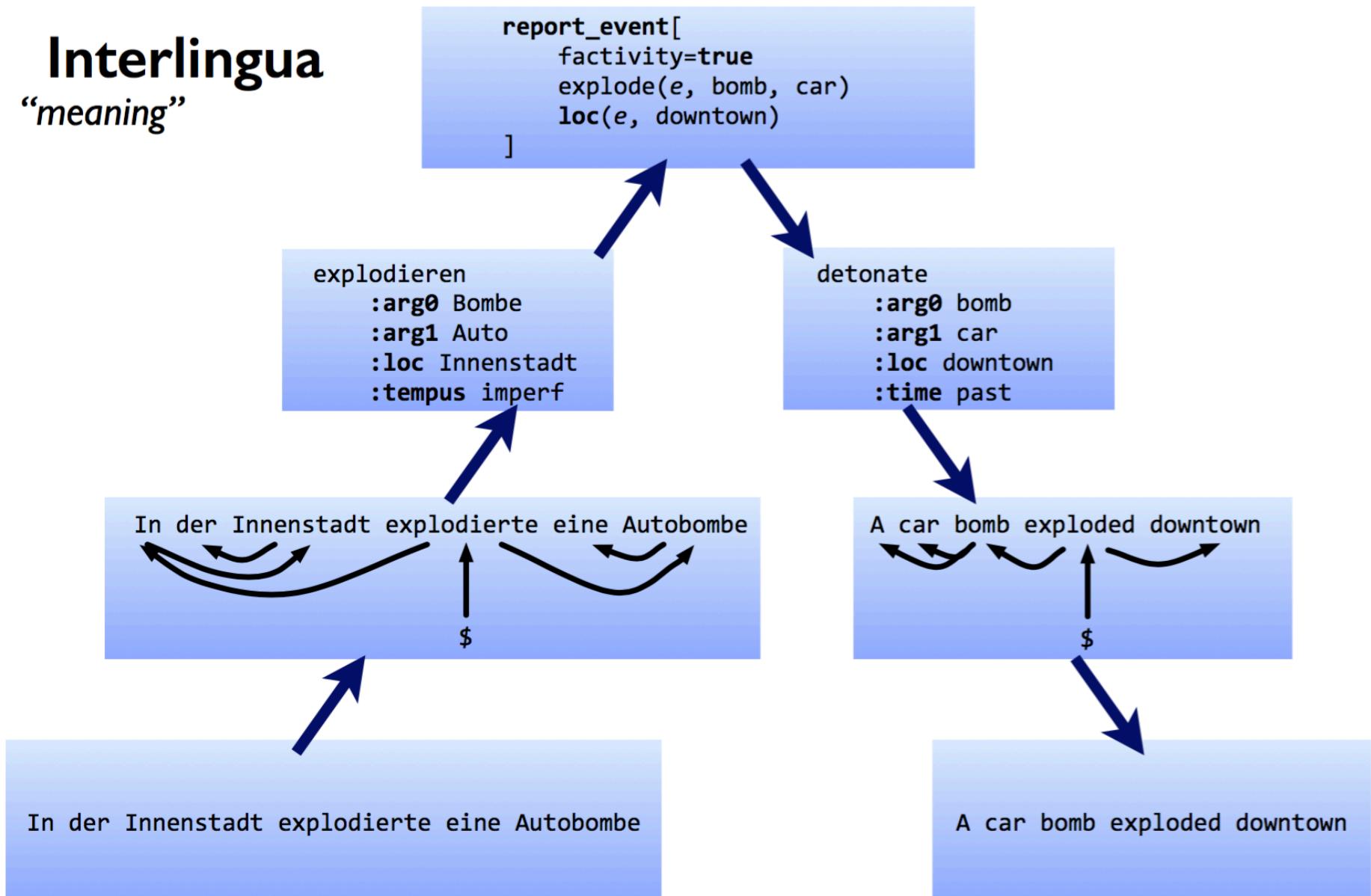


# Translation Modeling



# Translation Modeling

Interlingua  
“meaning”



# Translation Modeling

Interlingua  
“meaning”

```
report_event[  
    factivity=true  
    explode(e, bomb, car)  
    loc(e, downtown)  
]
```

```
explodieren  
:arg0 Bomb  
:arg1 car  
:loc Innenstadt  
:tempus imperf
```

```
detonate  
:arg0 bomb  
:arg1 car  
:loc downtown  
:time past
```

# Hidden

In der Innenstadt explodierte eine Autobombe

A car bomb exploded downtown

In der Innenstadt explodierte eine Autobombe

A car bomb exploded downtown

# Evaluation: many right answers

这个 机场 的 安全 工作 由 以色列 方面 负责 .

Israeli officials are responsible for airport security.

Israel is in charge of the security at this airport.

The security work for this airport is the responsibility of the Israel government.

Israeli side was in charge of the security of this airport.

Israel is responsible for the airport's security.

Israel is responsible for safety work at this airport.

Israel presides over the security of the airport.

Israel took charge of the airport security.

The safety of this airport is taken charge of by Israel.

This airport's security is the responsibility of the Israeli security officials.

# Evaluation: BLEU

N-Gram  
precision

$$p_n = \frac{\sum_{n\text{-gram} \in \text{hyp}} \text{count}_{clip}(n\text{-gram})}{\sum_{n\text{-gram} \in \text{hyp}} \text{count}(n\text{-gram})}$$

Bounded above  
by highest count  
of n-gram in any  
reference sentence

brevity  
penalty

$$B = \begin{cases} e^{(1 - |\text{ref}| / |\text{hyp}|)} & \text{if } |\text{ref}| > |\text{hyp}| \\ 1 & \text{otherwise} \end{cases}$$

Bleu score:  
brevity penalty,  
geometric  
mean of N-Gram  
precisions

$$\text{Bleu} = B \cdot \exp \left[ \frac{1}{N} \sum_{n=1}^N p_n \right]$$

# Plan

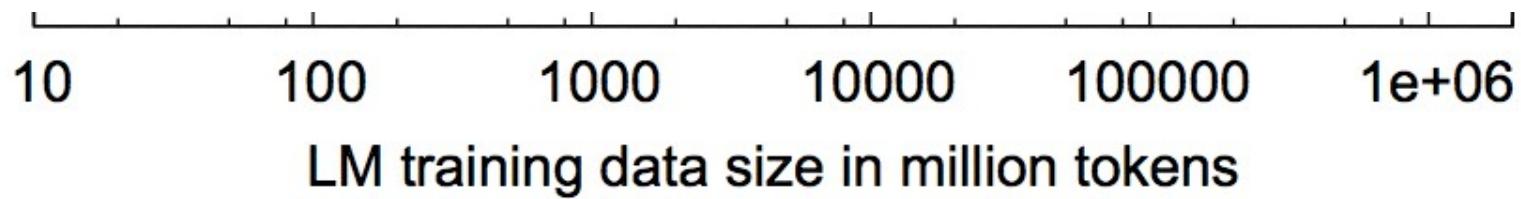
- Language models
- Word-based translation models
- Phrase-based translation models
- Decoding in PBSMT
- Discriminative modeling
- Neural machine translation

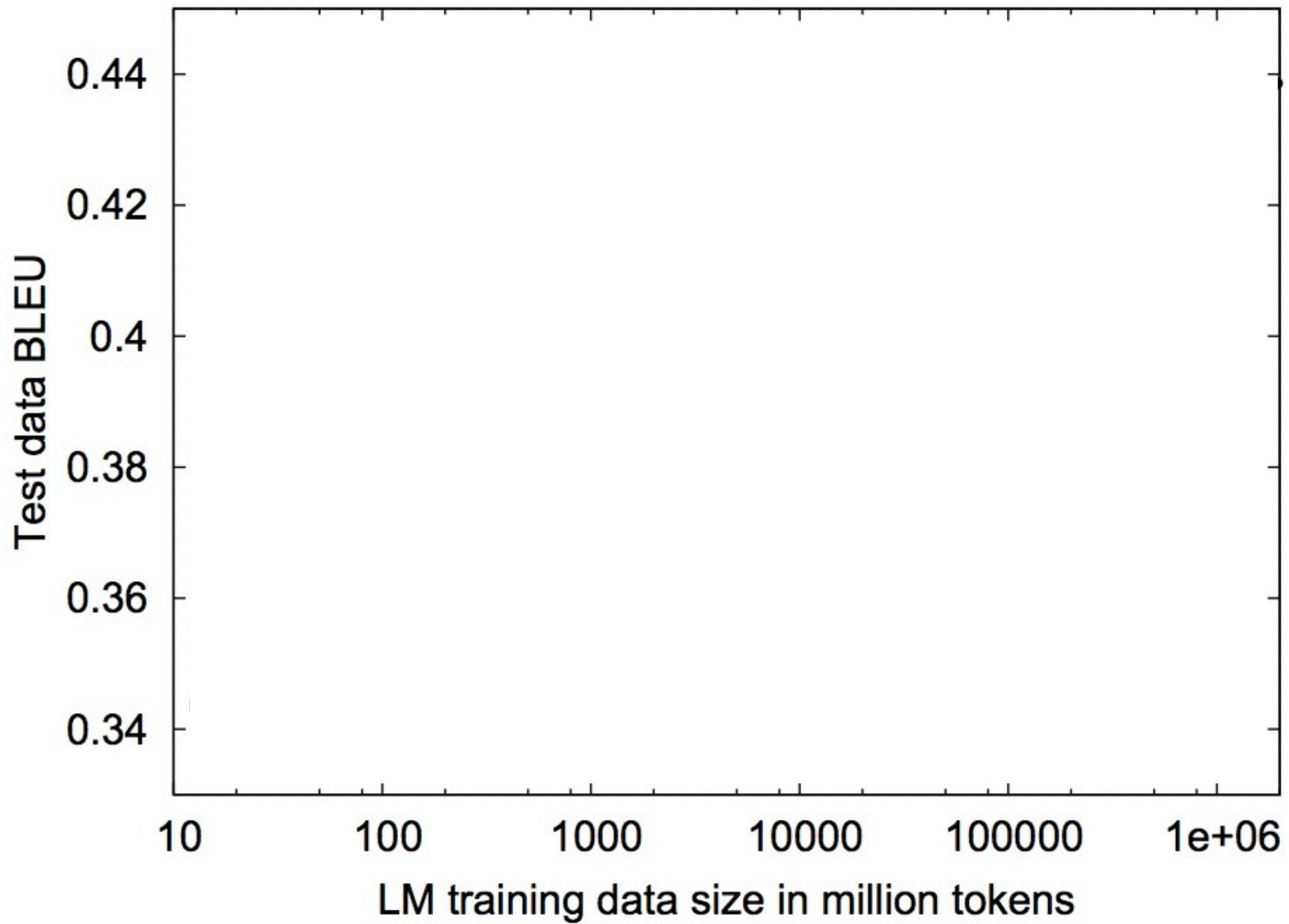
# Language Models (Lecture 5)

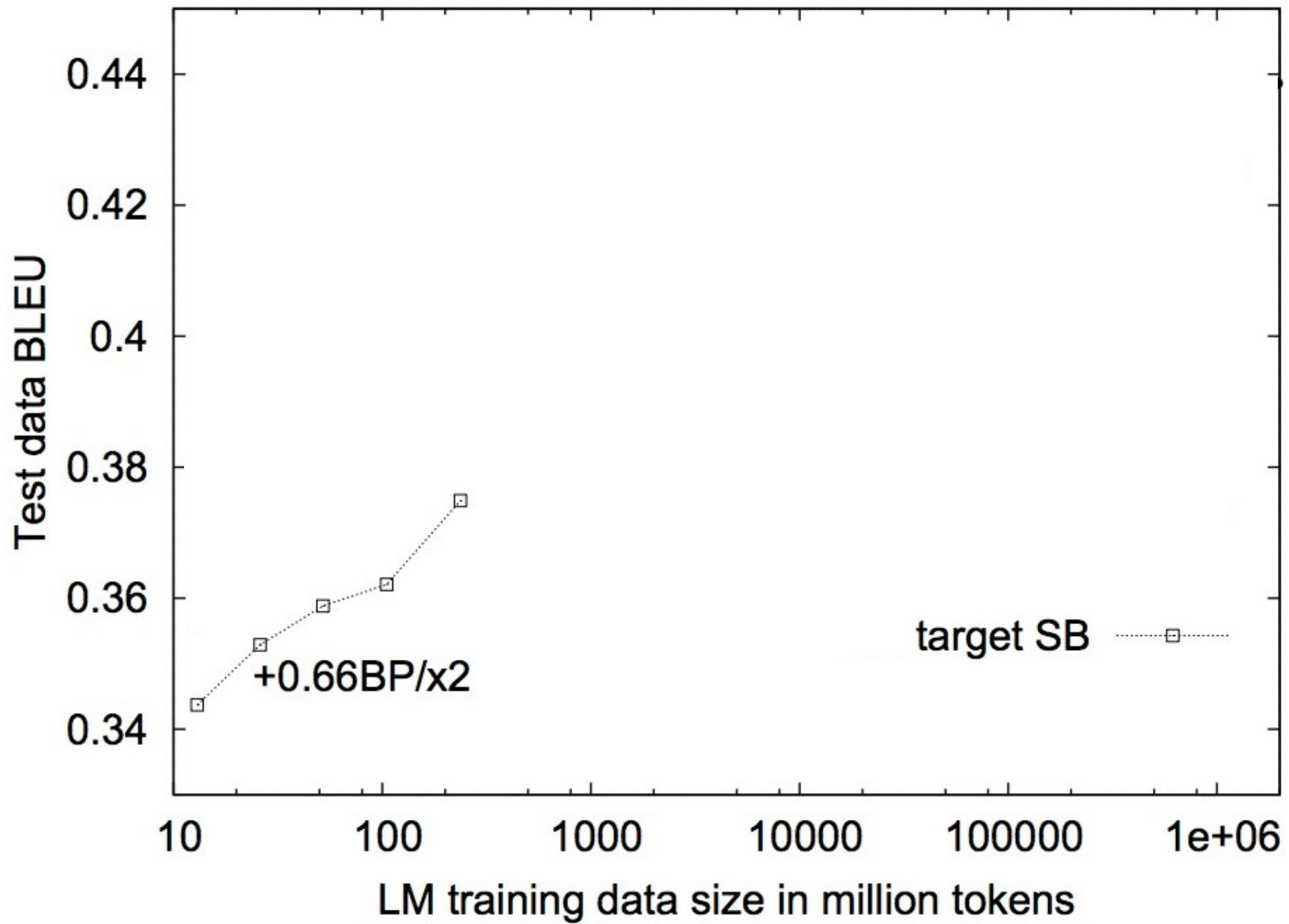
# Language Models

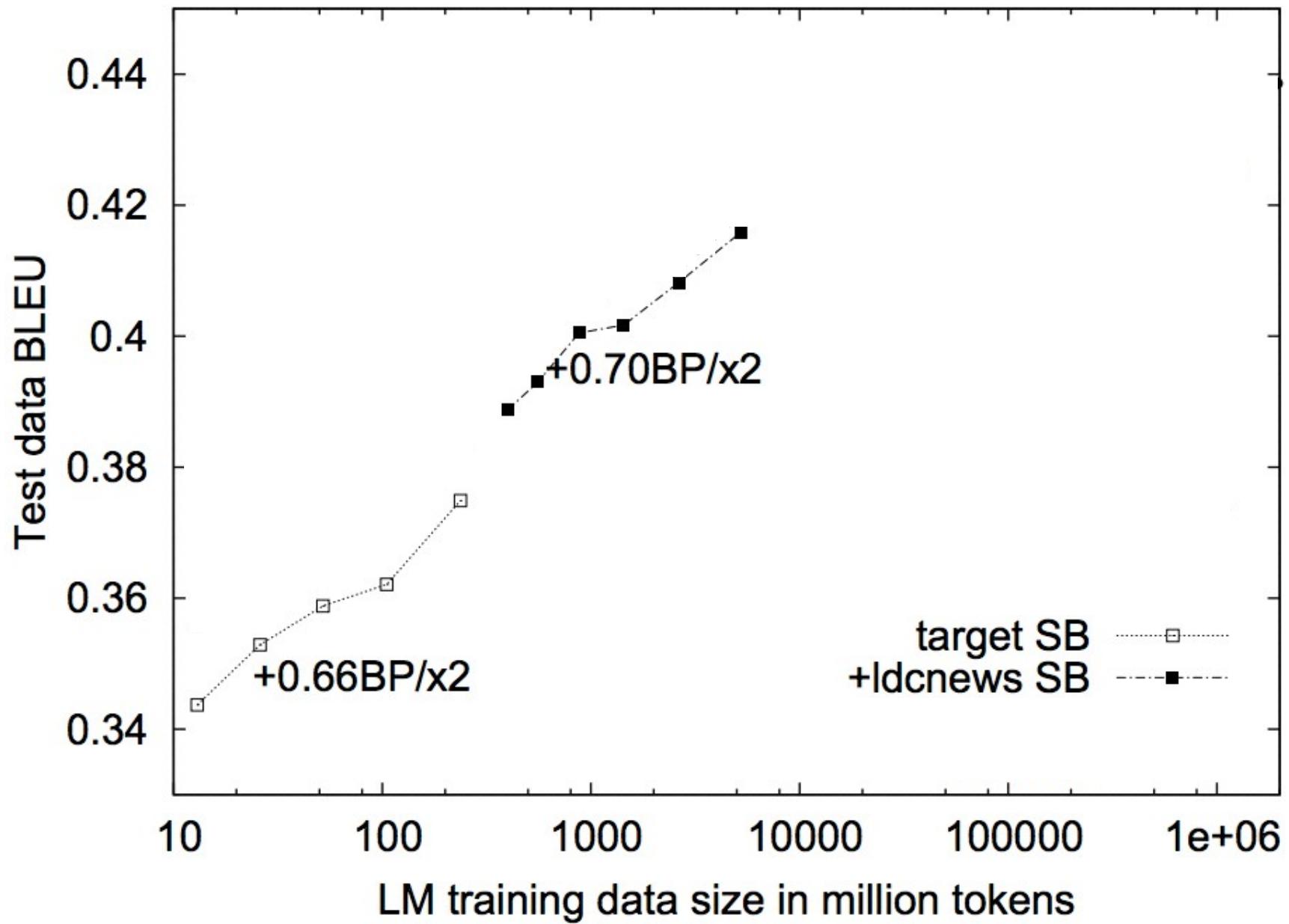
- Models of  $p(\text{sentence})$
- Later  
 $p(\text{translated sentence} \mid \text{input sentence})$
- Why do it this way?
  - Conditioning on more stuff makes modeling more complicated. That is:  $p(\text{sentence})$  is easier than  $p(\text{translated sentence} \mid \text{input sentence})$ .
  - Language models are arguably the most important models in statistical MT

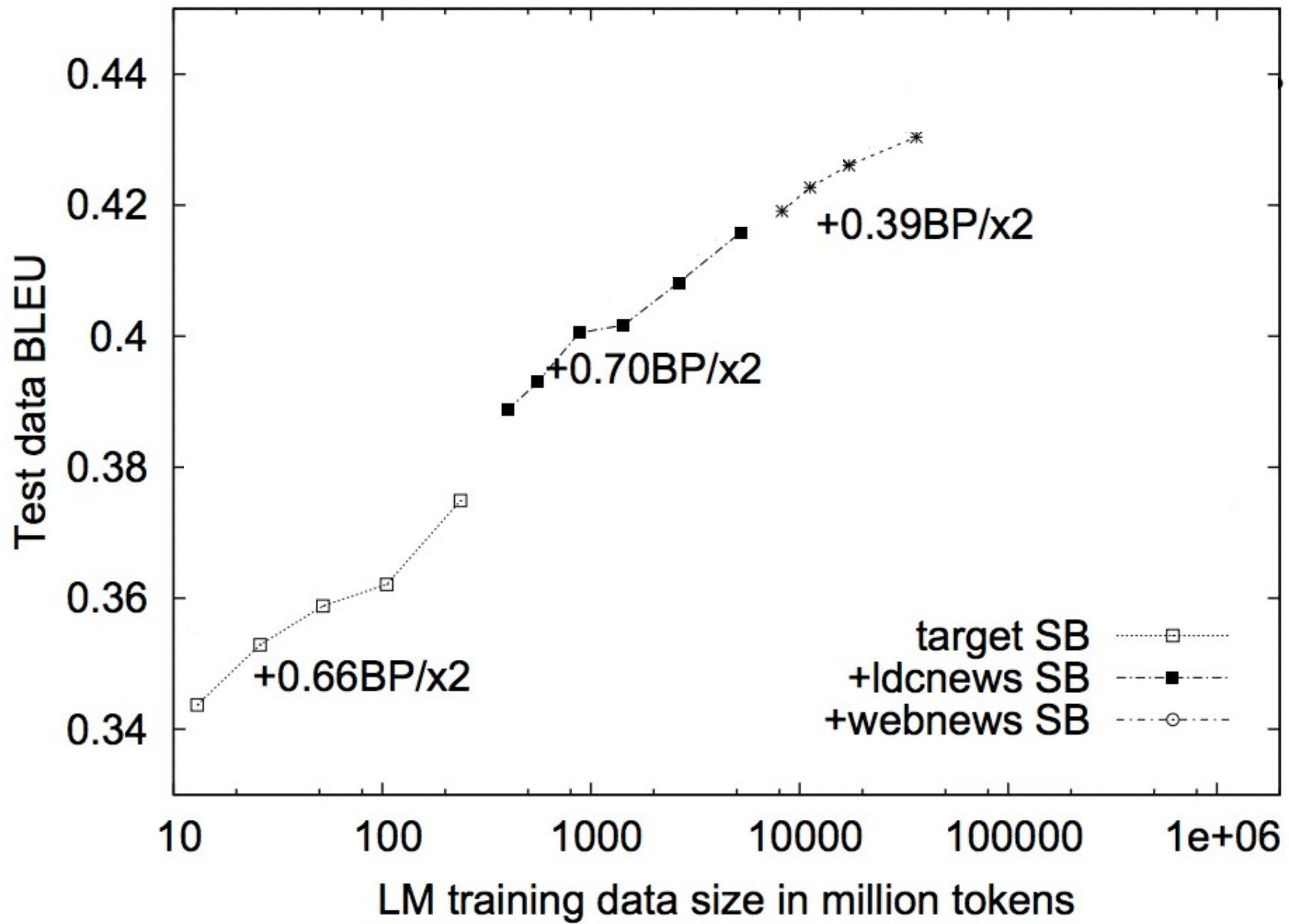
LM training data size in million tokens

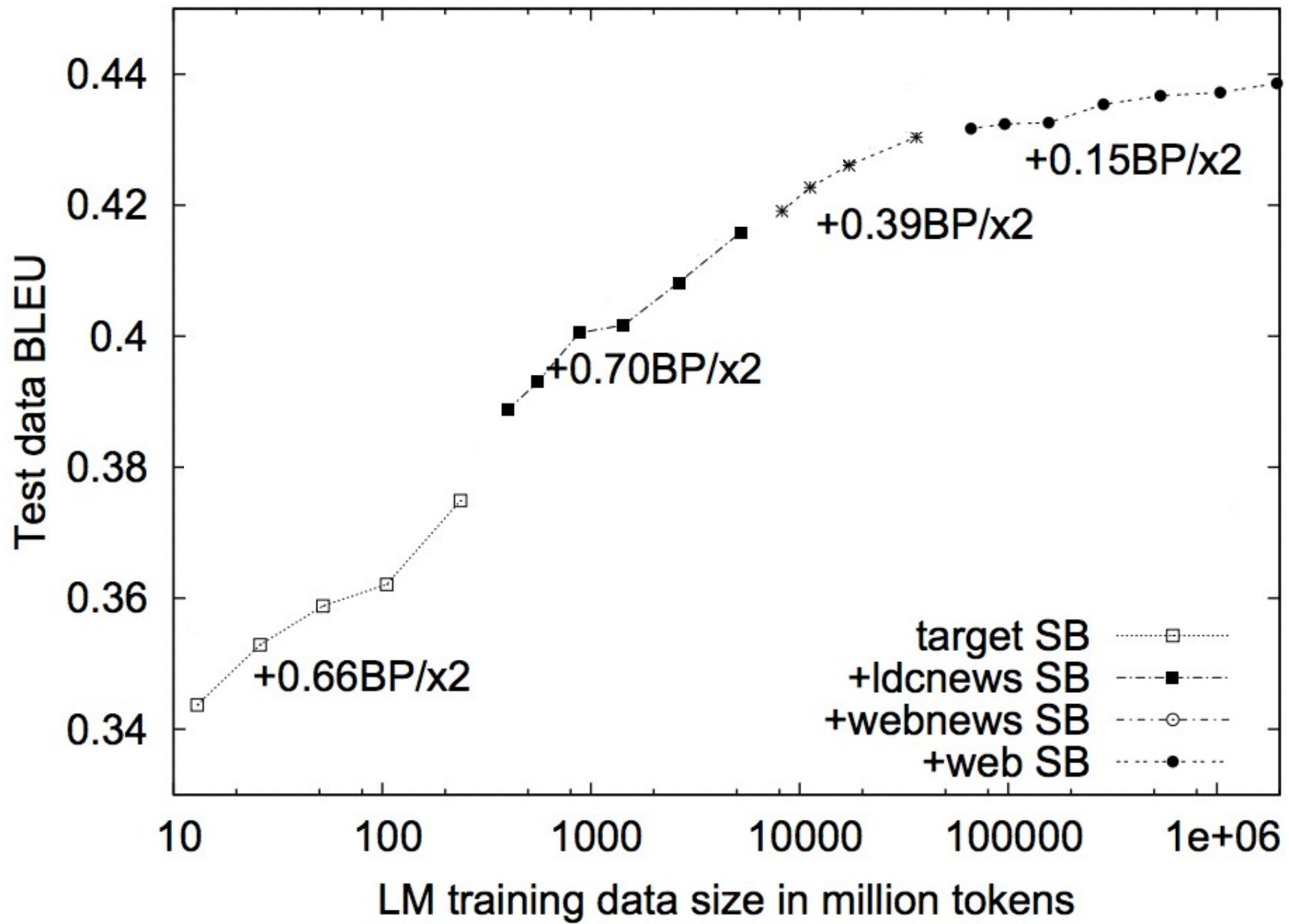












# Language Models Matter

- Language models play the role of ...
  - a judge of *grammaticality*
  - a judge of *semantic plausibility*
  - an enforcer of *stylistic consistency*

# Lexical (word-based) Translation Models

# Lexical Translation

- How do we translate a word? Look it up in the dictionary

*Haus* : *house, home, shell, household*

- Multiple translations
  - Different word senses, different registers, different inflections (?)
  - *house, home* are common
  - *shell* is specialized (the Haus of a snail is a shell)

# How common is each?

Translation	Count
house	5000
home	2000
shell	100
household	80

# MLE

$$\hat{p}_{\text{MLE}}(e \mid \text{Haus}) = \begin{cases} 0.696 & \text{if } e = \text{house} \\ 0.279 & \text{if } e = \text{home} \\ 0.014 & \text{if } e = \text{shell} \\ 0.011 & \text{if } e = \text{household} \\ 0 & \text{otherwise} \end{cases}$$

# Lexical Translation

- Goal: a model  $p(e | f, m)$
- where  $e$  and  $f$  are complete English and Foreign sentences

$$e = \langle e_1, e_2, \dots, e_m \rangle \quad f = \langle f_1, f_2, \dots, f_n \rangle$$

The diagram consists of two equations. The first equation shows  $e$  as a sequence of elements  $e_1, e_2, \dots, e_m$ . The second equation shows  $f$  as a sequence of elements  $f_1, f_2, \dots, f_n$ . Two blue arrows point upwards from the right side of each equation towards the corresponding sequence of elements.

# Lexical Translation

- Goal: a model  $p(\mathbf{e} \mid \mathbf{f}, m)$
- where  $\mathbf{e}$  and  $\mathbf{f}$  are complete English and Foreign sentences
- Lexical translation makes the following *assumptions*:
  - Each word in  $e_i$  in  $\mathbf{e}$  is generated from exactly one word in  $\mathbf{f}$
  - Thus, we have an *alignment*  $a_i$  that indicates which word  $e_i$  “came from”, specifically it came from  $f_{a_i}$ .
  - Given the alignments  $\mathbf{a}$ , translation decisions are conditionally independent of each other and depend *only* on the aligned source word  $f_{a_i}$ .

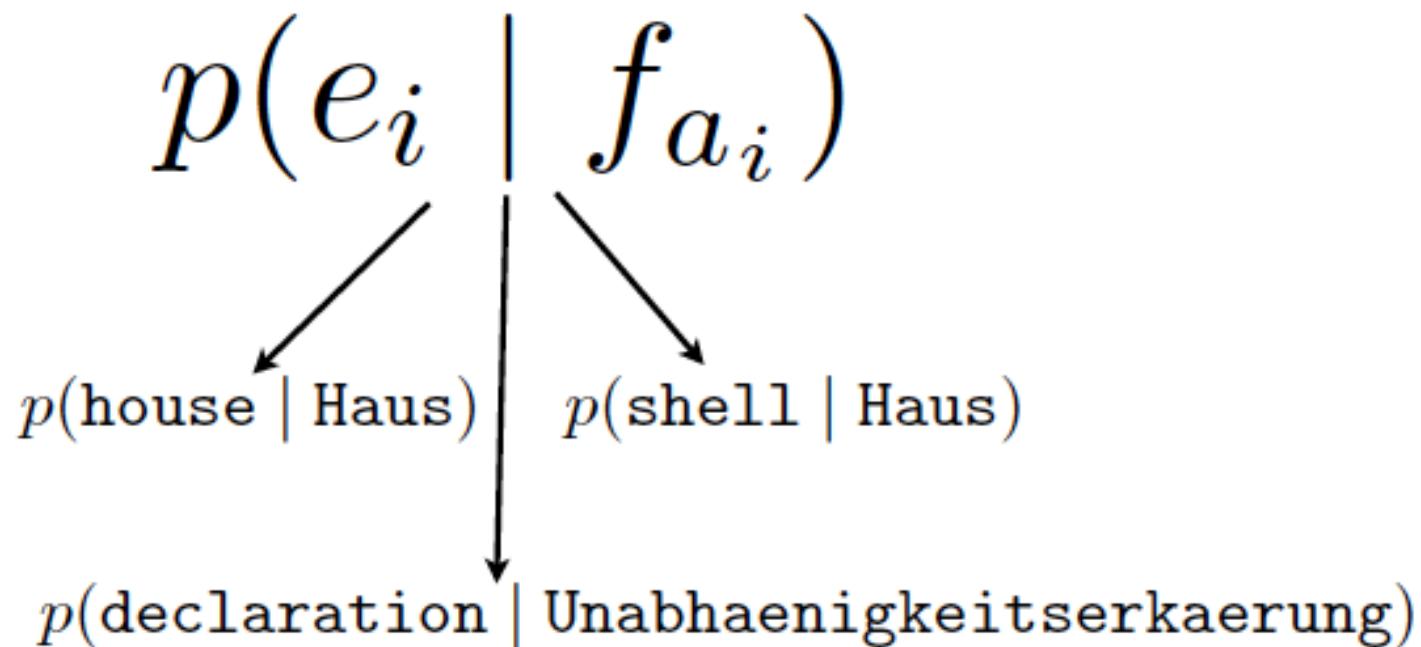
# Lexical Translation

- Putting our assumptions together, we have:

$$p(\mathbf{e} \mid \mathbf{f}, m) = \sum_{\mathbf{a} \in [0, n]^m} p(\mathbf{a} \mid \mathbf{f}, m) \times \prod_{i=1}^m p(e_i \mid f_{a_i})$$

Alignment  $\times$  Translation | Alignment

# Lexical Translation



Remember bigram models...

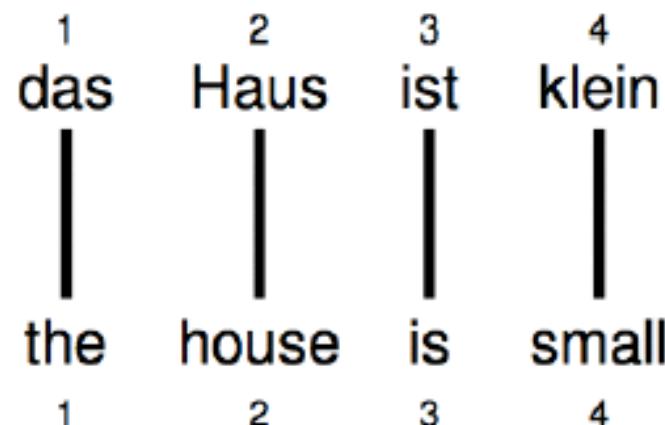
# Alignment

$$p(\mathbf{a} \mid \mathbf{f}, m)$$

Most of the action for the first 10 years of MT was here. Words weren't the problem, word *order* was hard.

# Alignment

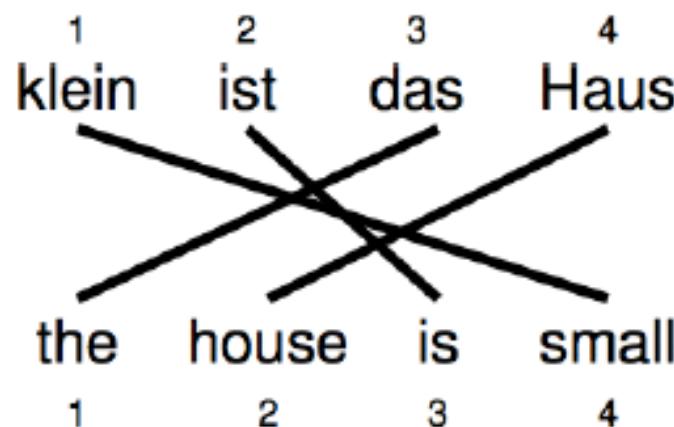
- Alignments can be visualized by drawing links between two sentences, and they are represented as vectors of positions:



$$\mathbf{a} = (1, 2, 3, 4)^\top$$

# Reordering

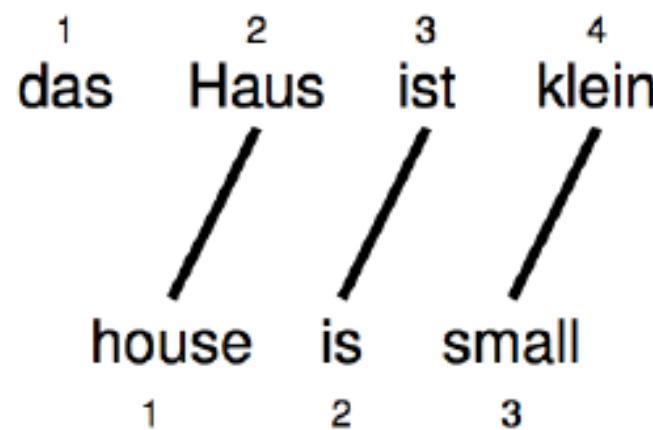
- Words may be reordered during translation.



$$\mathbf{a} = (3, 4, 2, 1)^\top$$

# Word Dropping

- A source word may not be translated at all

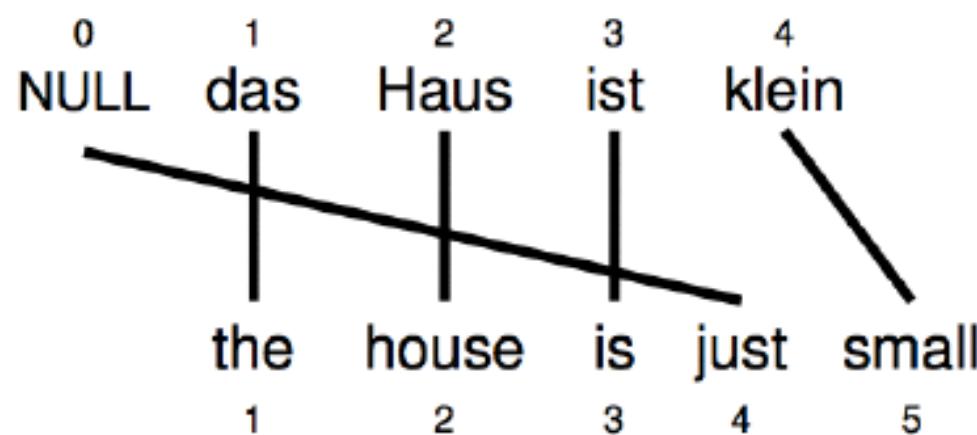


$$\mathbf{a} = (2, 3, 4)^\top$$

# Word Insertion

- Words may be inserted during translation  
English **just** does not have an equivalent

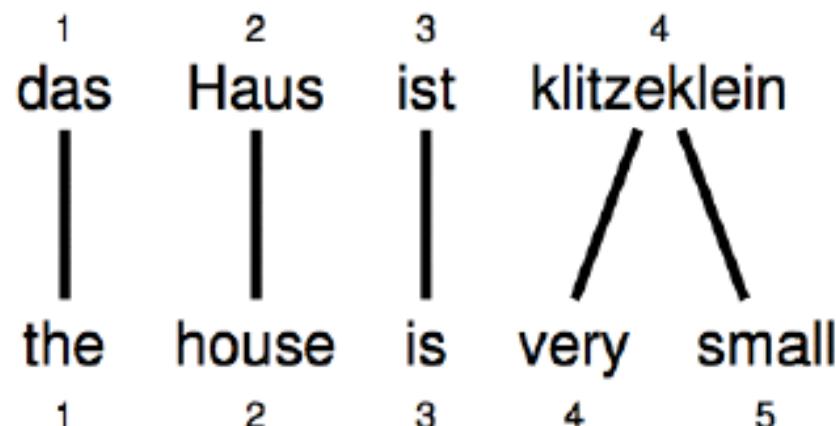
But it must be explained - we typically assume every source sentence contains a NULL token



$$\mathbf{a} = (1, 2, 3, 0, 4)^\top$$

# One-to-many Translation

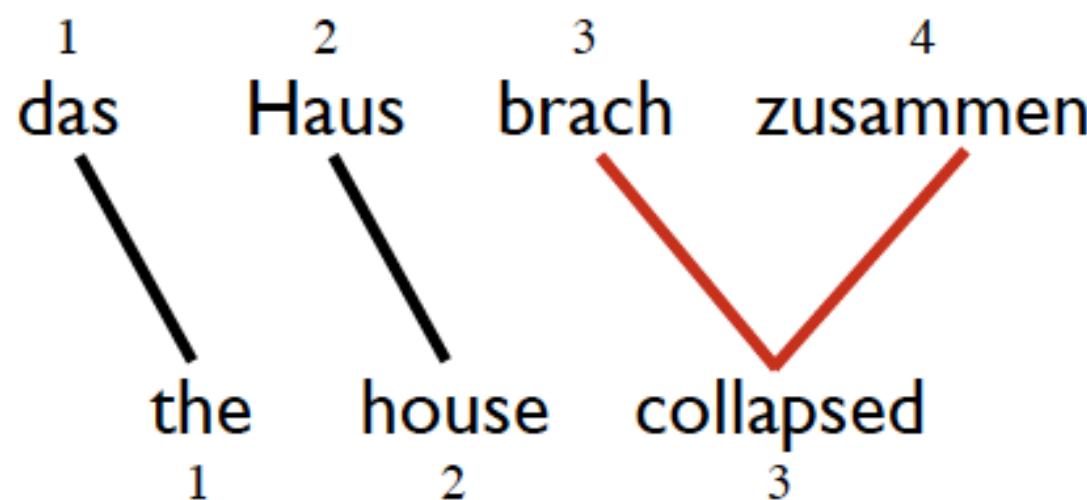
- A source word may translate into **more than one** target word



$$\mathbf{a} = (1, 2, 3, 4, 4)^\top$$

# Many-to-one Translation

- More than one source word may not translate as a unit in lexical translation



a = ???

# IBM Model I

- Simplest possible lexical translation model
- Additional assumptions
  - The  $m$  alignment decisions are independent
  - The alignment distribution for each  $a_i$  is uniform over all source words and NULL

for each  $i \in [1, 2, \dots, m]$

$$a_i \sim \text{Uniform}(0, 1, 2, \dots, n)$$

$$e_i \sim \text{Categorical}(\theta_{f_{a_i}})$$

# IBM Model I

for each  $i \in [1, 2, \dots, m]$

$$a_i \sim \text{Uniform}(0, 1, 2, \dots, n)$$

$$e_i \sim \text{Categorical}(\theta_{f_{a_i}})$$

$$p(\mathbf{e}, \mathbf{a} \mid \mathbf{f}, m) = \prod_{i=1}^m$$

# IBM Model I

for each  $i \in [1, 2, \dots, m]$

$$a_i \sim \text{Uniform}(0, 1, 2, \dots, n)$$

$$e_i \sim \text{Categorical}(\theta_{f_{a_i}})$$

$$p(\mathbf{e}, \mathbf{a} \mid \mathbf{f}, m) = \prod_{i=1}^m \frac{1}{1 + n}$$

# IBM Model I

for each  $i \in [1, 2, \dots, m]$

$$a_i \sim \text{Uniform}(0, 1, 2, \dots, n)$$

$$e_i \sim \text{Categorical}(\theta_{f_{a_i}})$$

$$p(\mathbf{e}, \mathbf{a} \mid \mathbf{f}, m) = \prod_{i=1}^m \frac{1}{1+n} p(e_i \mid f_{a_i})$$

# Example

0      1      2      3      4  
NULL    das    Haus    ist    klein

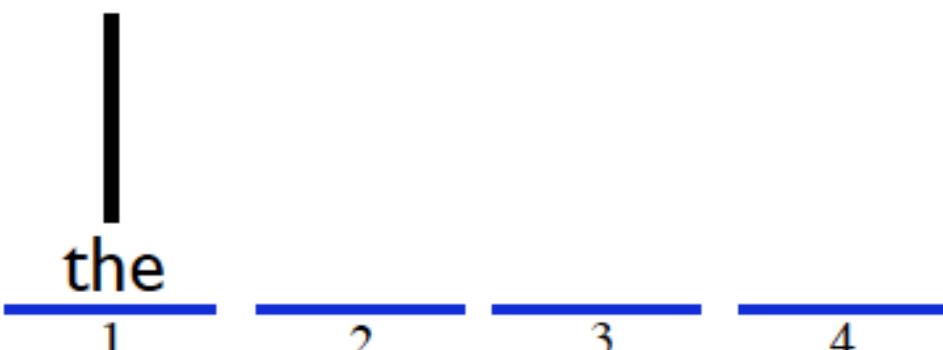
The diagram illustrates a sequence of words with corresponding indices above them. The words are: NULL, das, Haus, ist, and klein. Above each word is its index: 0, 1, 2, 3, and 4 respectively. Below the sequence, there are four horizontal blue bars, each aligned with one of the words. Below each bar is its index: 1, 2, 3, and 4.

Start with a foreign sentence and a target length.

# Example

0      1      2      3      4  
NULL    das    Haus    ist    klein  
  
1      2      3      4

# Example

0      1      2      3      4  
NULL    das    Haus    ist    klein  
  
the  
1      2      3      4

# Example

0	1	2	3	4
NULL	das	Haus	ist	klein
	the			
	1	2	3	4

The diagram illustrates word embeddings for the German sentence "das Haus ist klein". The words are indexed from 0 to 4. The word "the" is mapped to index 1, "das" to index 2, "Haus" to index 3, "ist" to index 4, and "klein" to index 0. Vertical lines connect the words to their respective indices. Below the indices, horizontal blue lines indicate the embedding dimensions for each word.

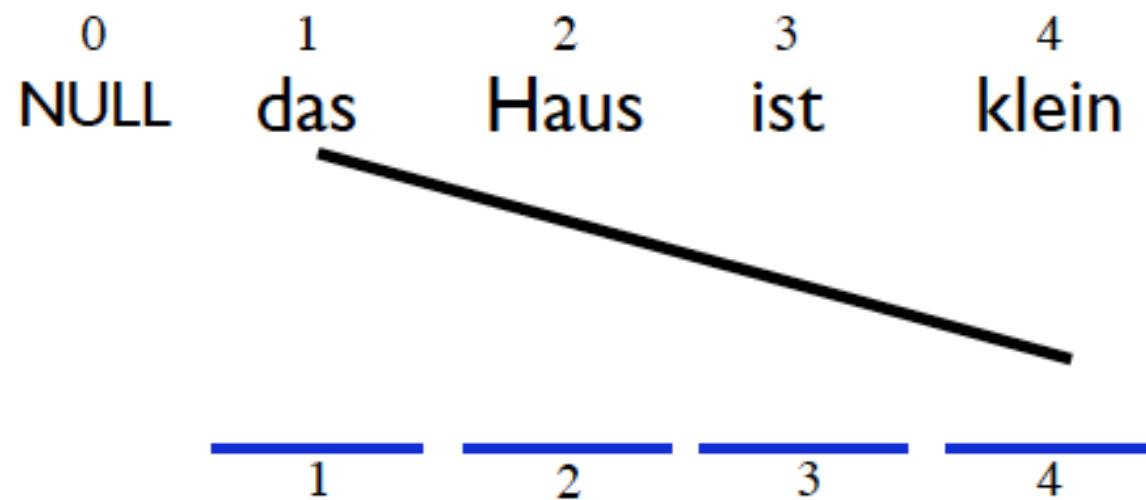
# Example

0	1	2	3	4
NULL	das	Haus	ist	klein
	<u>the</u>	<u>house</u>	<u>3</u>	<u>4</u>

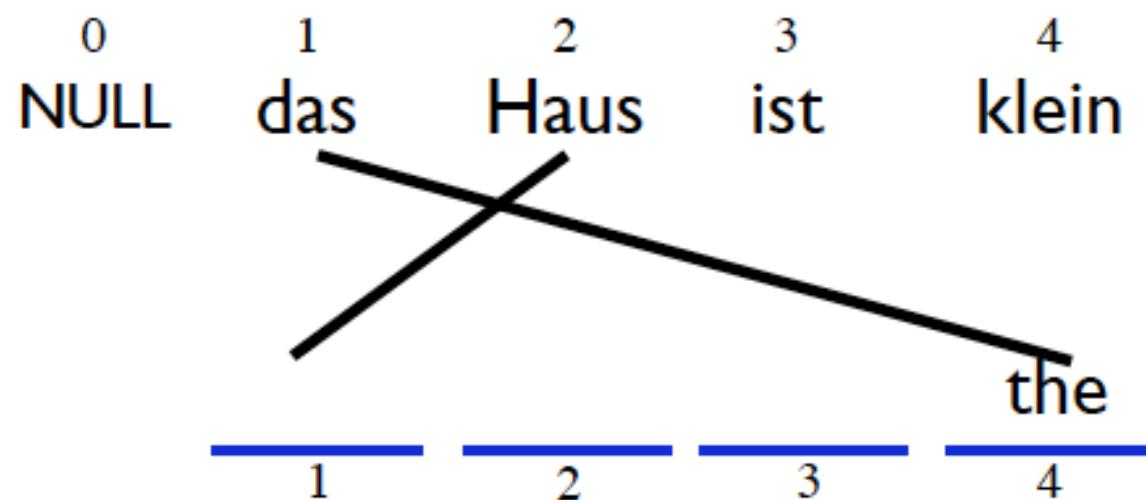
# Example

0	1	2	3	4
NULL	das	Haus	ist	klein
				
	<u>the</u>	<u>house</u>	<u>is</u>	<u>small</u>
	1	2	3	4

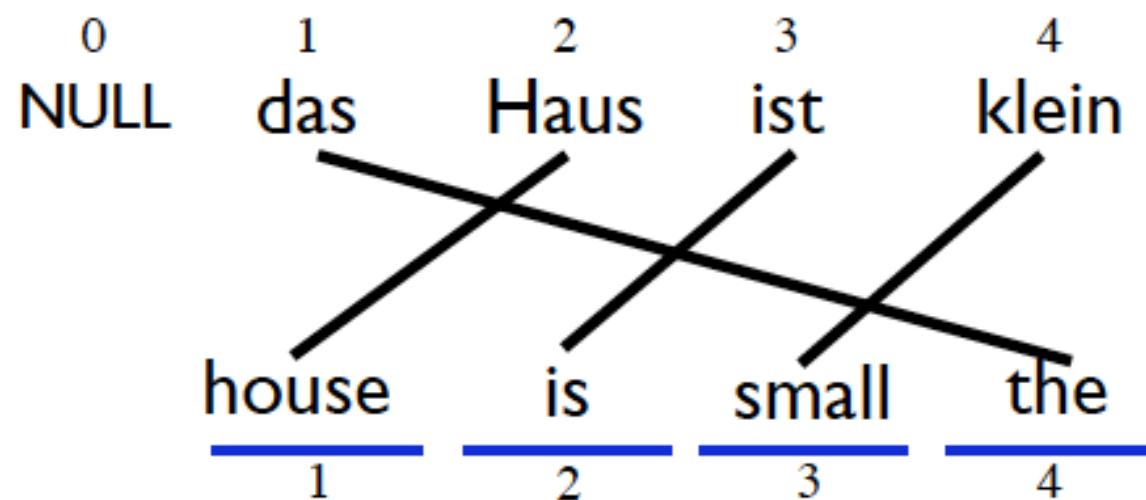
# Example



# Example



# Example



# Finding the Viterbi Alignment

$$\begin{aligned}\mathbf{a}^* &= \arg \max_{\mathbf{a} \in [0,1,\dots,n]^m} p(\mathbf{a} \mid \mathbf{e}, \mathbf{f}) \\ &= \arg \max_{\mathbf{a} \in [0,1,\dots,n]^m} \frac{p(\mathbf{e}, \mathbf{a} \mid \mathbf{f})}{\sum_{\mathbf{a}'} p(\mathbf{e}, \mathbf{a}' \mid \mathbf{f})} \\ &= \arg \max_{\mathbf{a} \in [0,1,\dots,n]^m} p(\mathbf{e}, \mathbf{a} \mid \mathbf{f})\end{aligned}$$

$$\begin{aligned}a_i^* &= \arg \max_{a_i=0}^n \frac{1}{1+n} p(e_i \mid f_{a_i}) \\ &= \arg \max_{a_i=0}^n p(e_i \mid f_{a_i})\end{aligned}$$

# Historical Note

The **Viterbi algorithm** is a dynamic programming algorithm for finding the most likely sequence of hidden states – called the **Viterbi path** – that results in a sequence of observed events, especially in the context of Markov information sources and hidden Markov models.

*Andrew Viterbi*

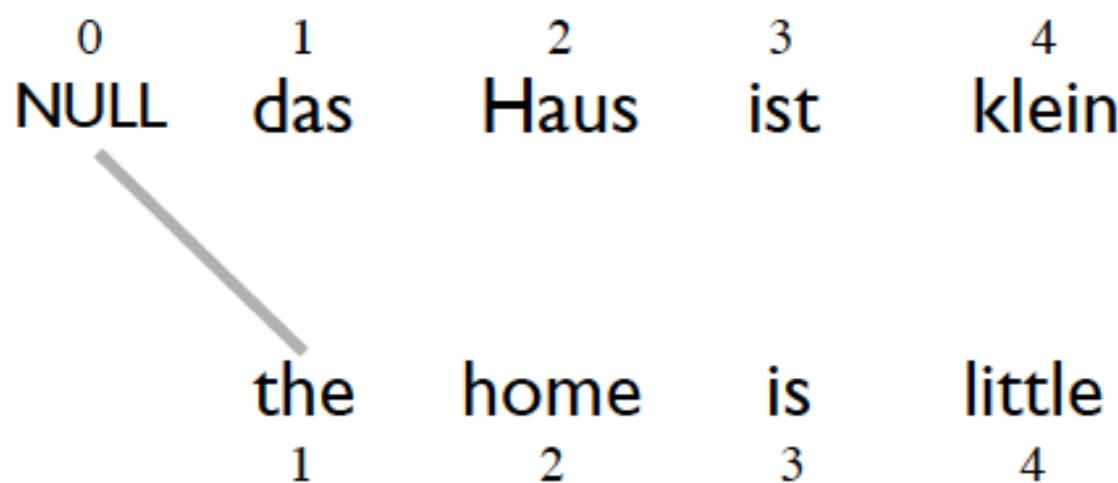
*Professor at USC*

*co-founder of Qualcomm*

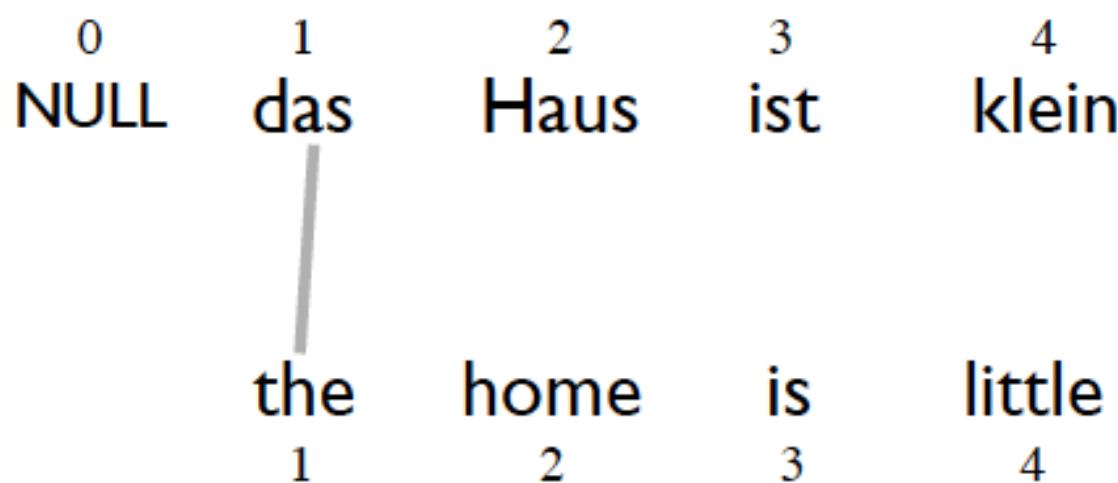
*classmates with Fred Jelinek*



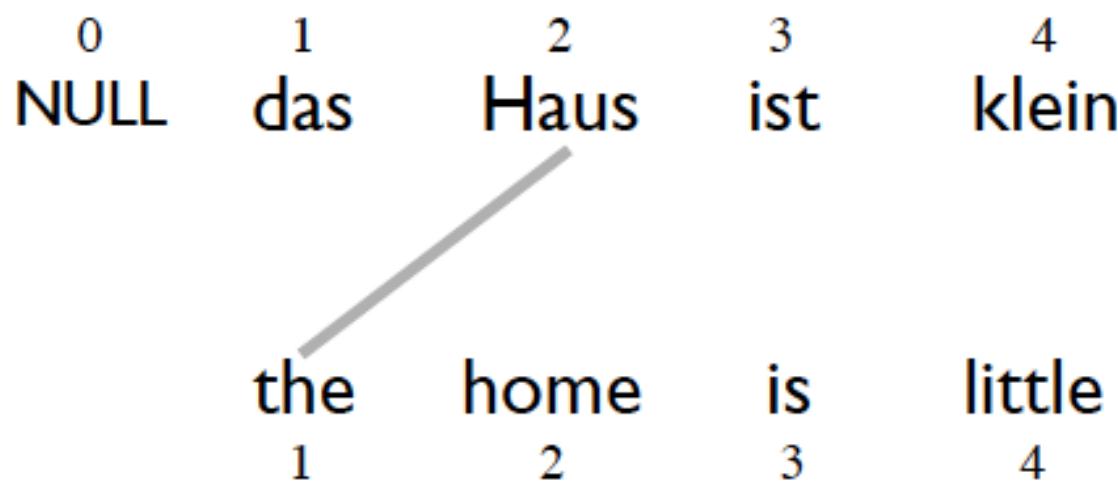
# Finding the Viterbi Alignment



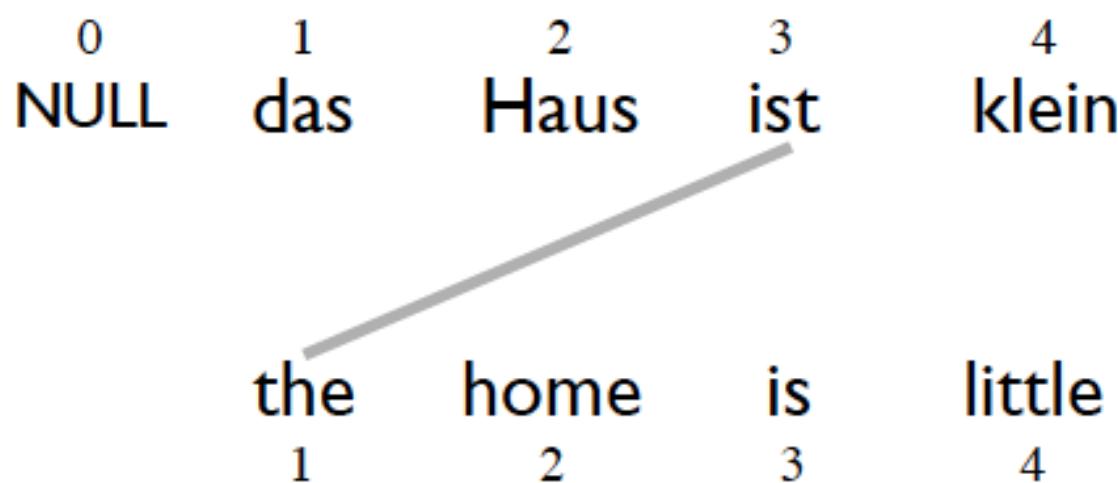
# Finding the Viterbi Alignment



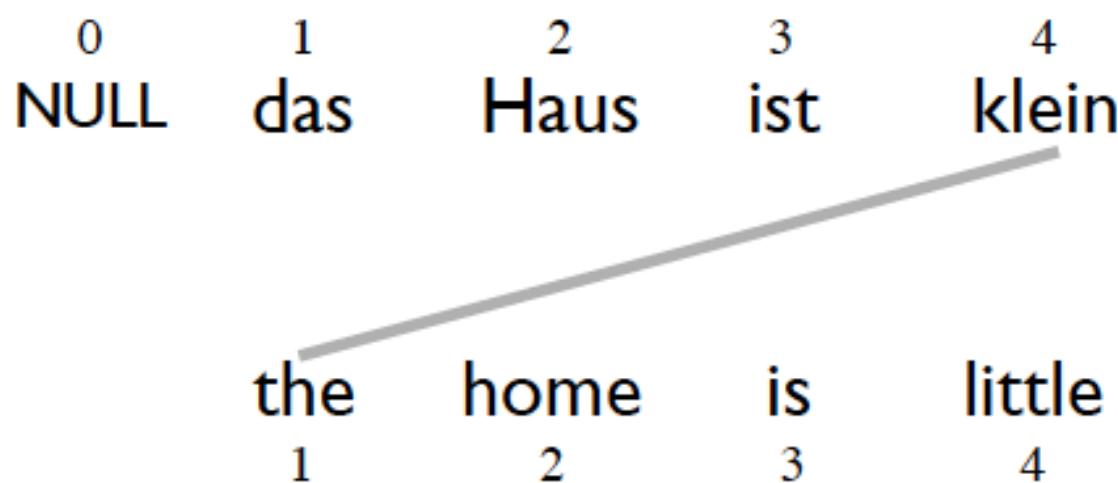
# Finding the Viterbi Alignment



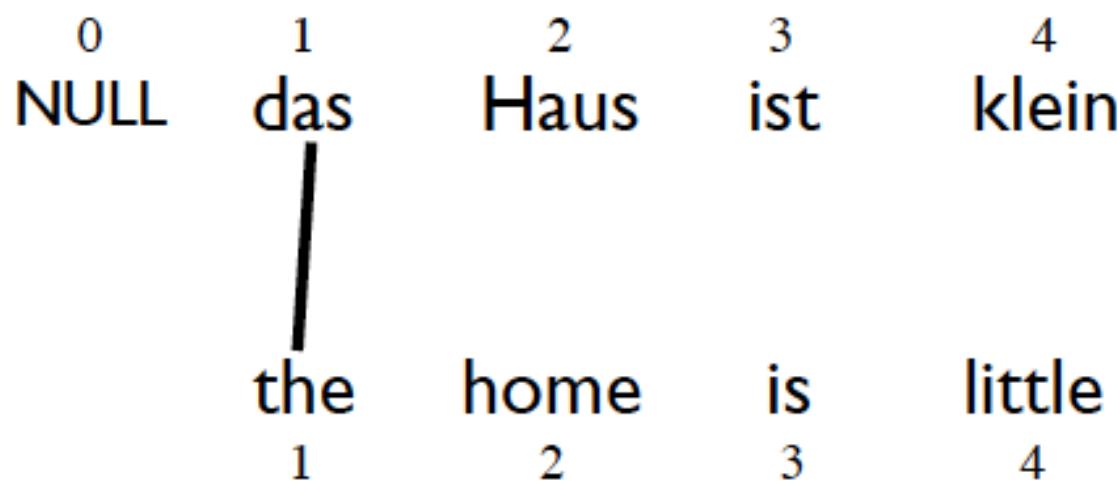
# Finding the Viterbi Alignment



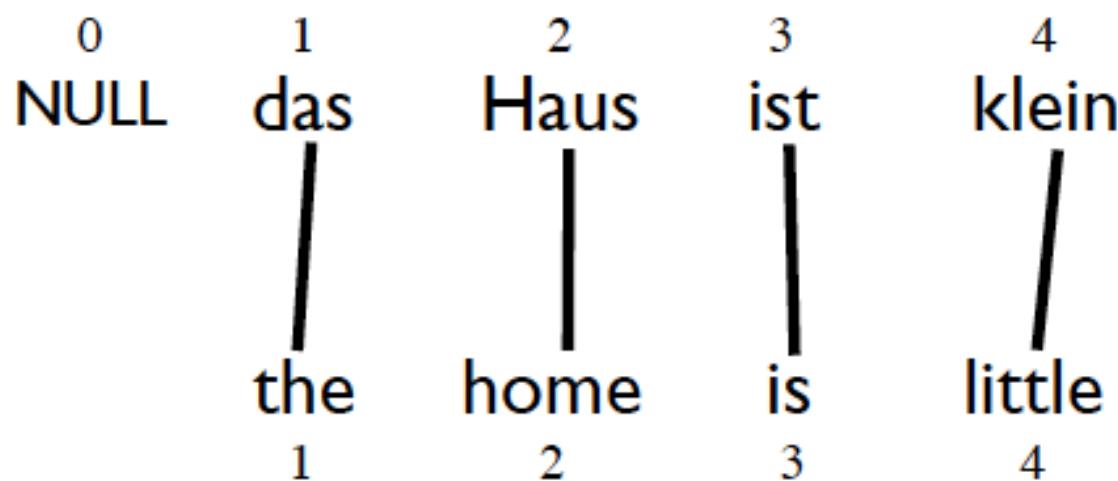
# Finding the Viterbi Alignment



# Finding the Viterbi Alignment

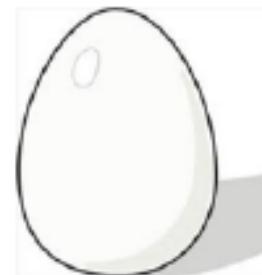
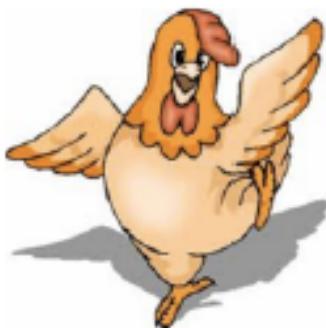


# Finding the Viterbi Alignment



# Learning Lexical Translation Models

- How do we learn the parameters  $p(e | f)$
- “Chicken and egg” problem
  - If we had the alignments, we could estimate the parameters (MLE)
  - If we had parameters, we could find the most likely alignments



# EM Algorithm

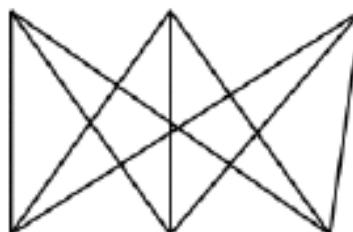
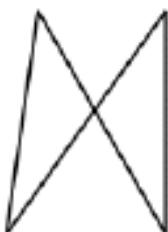
- pick some random (or uniform) parameters
- Repeat until you get bored (~ 5 iterations for lexical translation models)
  - using your current parameters, compute “expected” alignments for every target word token in the training data

$$p(a_i \mid e, f)$$

- keep track of the expected number of times  $f$  translates into  $e$  throughout the whole corpus
- keep track of the expected number of times that  $f$  is used as the source of any translation
- use these expected counts as if they were “real” counts in the standard MLE equation

# EM for Model I

... la maison ... la maison blue ... la fleur ...

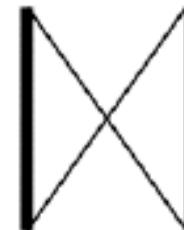
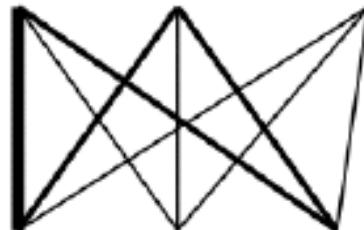


... the house ... the blue house ... the flower ...

- Initial step: all alignments equally likely
- Model learns that, e.g., la is often aligned with the

# EM for Model I

... la maison ... la maison blue ... la fleur ...

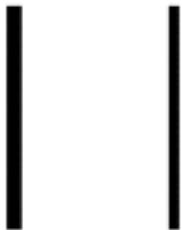


... the house ... the blue house ... the flower ...

- After one iteration
- Alignments, e.g., between **la** and **the** are more likely

# EM for Model I

... la maison ... la maison bleu ... la fleur ...



... the house ... the blue house ... the flower ...

- After another iteration
- It becomes apparent that alignments, e.g., between *fleur* and *flower* are more likely (pigeon hole principle)

# EM for Model I

... la maison ... la maison bleu ... la fleur ...  
  
... the house ... the blue house ... the flower ...

- Convergence
- Inherent hidden structure revealed by EM

# EM for Model I

... la maison ... la maison bleu ... la fleur ...



... the house ... the blue house ... the flower ...



$$p(\text{la}|\text{the}) = 0.453$$

$$p(\text{le}|\text{the}) = 0.334$$

$$p(\text{maison}|\text{house}) = 0.876$$

$$p(\text{bleu}|\text{blue}) = 0.563$$

...

- Parameter estimation from the aligned corpus

# Convergence

das Haus  
the house

das Buch  
the book

ein Buch  
a book

$e$	$f$	initial	1st it.	2nd it.	3rd it.	...	final
the	das	0.25	0.5	0.6364	0.7479	...	1
book	das	0.25	0.25	0.1818	0.1208	...	0
house	das	0.25	0.25	0.1818	0.1313	...	0
the	buch	0.25	0.25	0.1818	0.1208	...	0
book	buch	0.25	0.5	0.6364	0.7479	...	1
a	buch	0.25	0.25	0.1818	0.1313	...	0
book	ein	0.25	0.5	0.4286	0.3466	...	0
a	ein	0.25	0.5	0.5714	0.6534	...	1
the	haus	0.25	0.5	0.4286	0.3466	...	0
house	haus	0.25	0.5	0.5714	0.6534	...	1

# Can we do better?

$$p(\text{Translation}) = \sum_{\text{Alignment}} p(\text{Alignment}, \text{Translation})$$

$$= \sum_{\text{Alignment}} p(\text{Alignment}) \times p(\text{Translation} | \text{Alignment})$$

$$p(\mathbf{e} | \mathbf{f}, m) = \sum_{\mathbf{a} \in [0, n]^m} p(\mathbf{a} | \mathbf{f}, m) \times \prod_{i=1}^m p(e_i | f_{a_i})$$

$$p(\mathbf{e} \mid \mathbf{f}, m) = \sum_{\mathbf{a} \in [0,n]^m} p(\mathbf{a} \mid \mathbf{f}, m) \times \prod_{i=1}^m p(e_i \mid f_{a_i}) \\ \prod_{i=1}^m p(e_i \mid f_{a_i}, f_{a_{i-1}}) \\ \prod_{i=1}^m p(e_i \mid f_{a_i}, f_{a_i-1}) \\ \prod_{i=1}^m p(e_i \mid f_{a_i}, e_{i-1})$$

$$p(\mathbf{e} \mid \mathbf{f}, m) = \sum_{\mathbf{a} \in [0,n]^m} p(\mathbf{a} \mid \mathbf{f}, m) \times \prod_{i=1}^m p(e_i \mid f_{a_i})$$

$$= \sum_{\mathbf{a} \in [0,n]^m} \underbrace{\prod_{i=1}^m \frac{1}{1+n}}_{p(\mathbf{a}|\mathbf{f},m)} \times \prod_{i=1}^m p(e_i \mid f_{a_i})$$

$$= \sum_{\mathbf{a} \in [0,n]^m} \prod_{i=1}^m \frac{1}{1+n} p(e_i \mid f_{a_i})$$

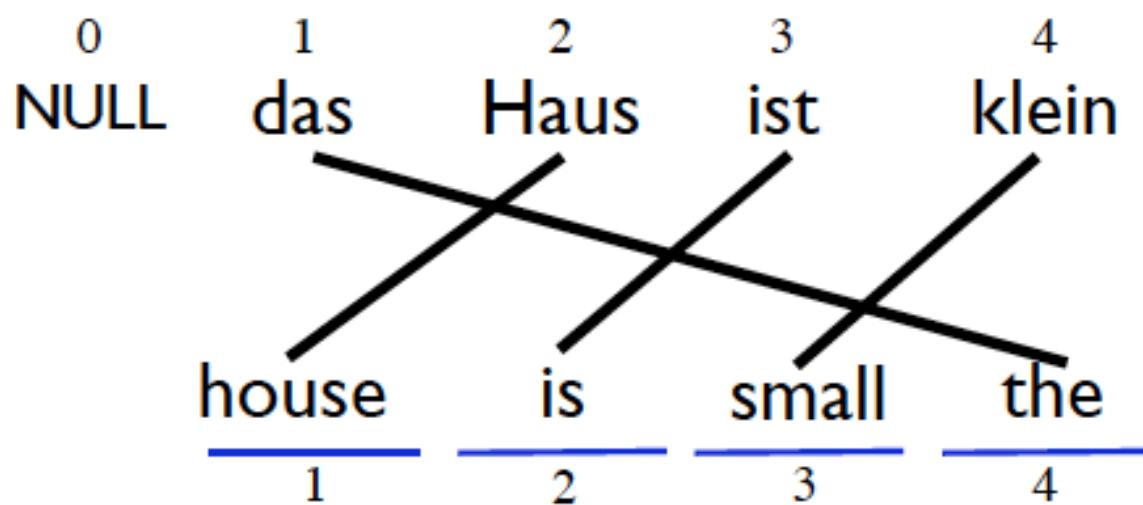
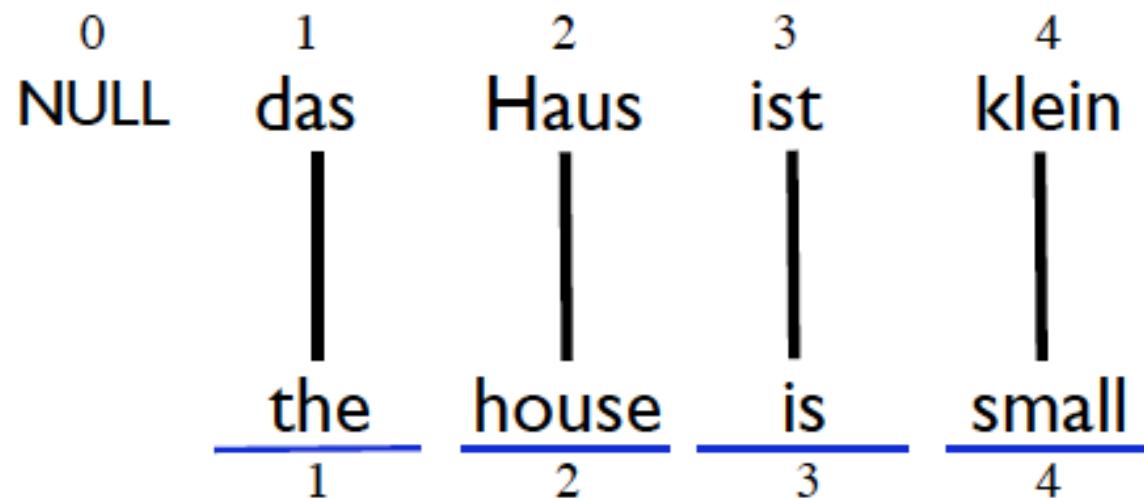
$$= \sum_{\mathbf{a} \in [0,n]^m} \prod_{i=1}^m p(a_i) \times p(e_i \mid f_{a_i})$$

$$p(\mathbf{e} \mid \mathbf{f}, m) = \sum_{\mathbf{a} \in [0, n]^m} p(\mathbf{a} \mid \mathbf{f}, m) \times \prod_{i=1}^m p(e_i \mid f_{a_i})$$

$$= \sum_{\mathbf{a} \in [0, n]^m} \underbrace{\prod_{i=1}^m \frac{1}{1+n}}_{p(\mathbf{a} \mid \mathbf{f}, m)} \times \prod_{i=1}^m p(e_i \mid f_{a_i})$$

**Can we do something better here?**

$$= \sum_{\mathbf{a} \in [0, n]^m} \prod_{i=1}^m p(a_i) \times p(e_i \mid f_{a_i})$$

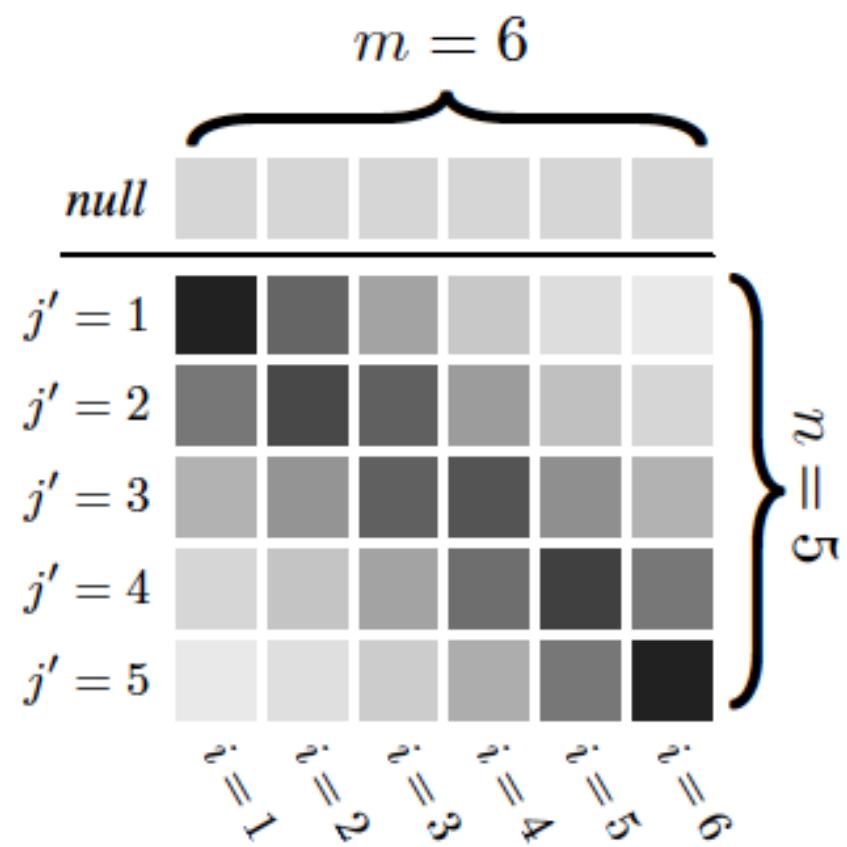


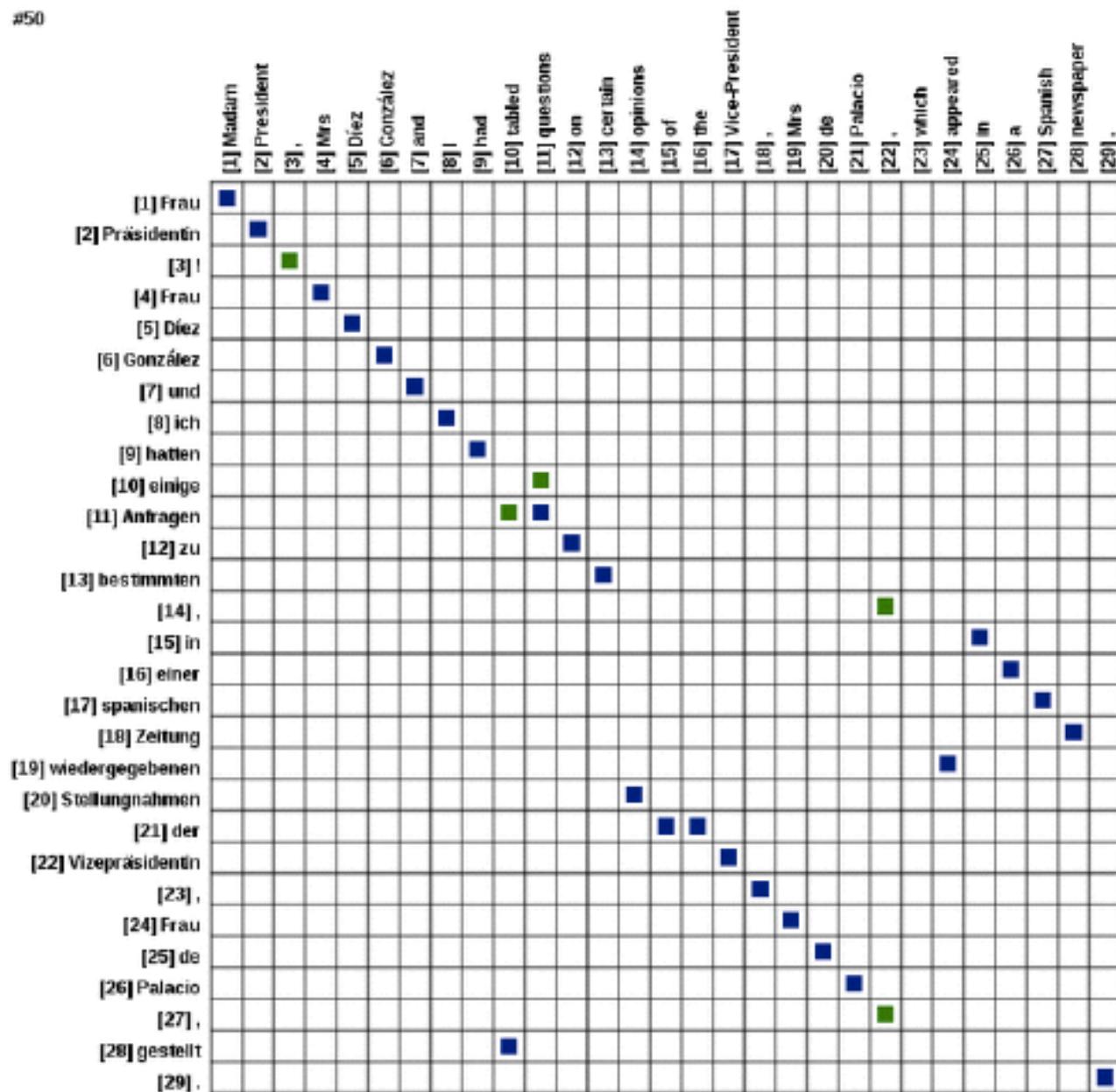
$$p(\mathbf{e} \mid \mathbf{f}, m) = \sum_{\mathbf{a} \in [0, n]^m} \prod_{i=1}^m p(a_i) \times p(e_i \mid f_{a_i})$$

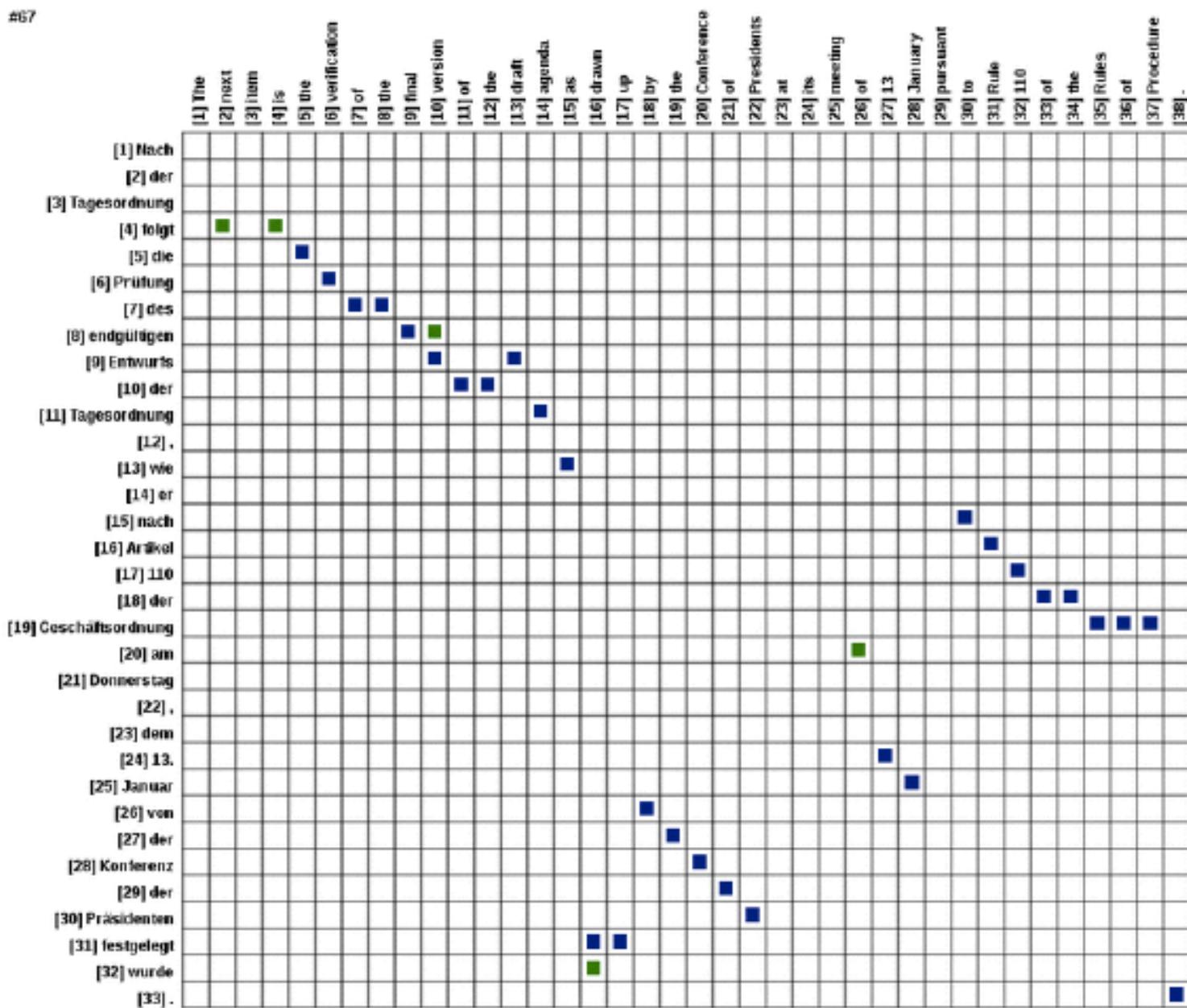
$$\text{Model 2} = \sum_{\mathbf{a} \in [0, n]^m} \prod_{i=1}^m p(a_i \mid i, m, n) \times p(e_i \mid f_{a_i})$$

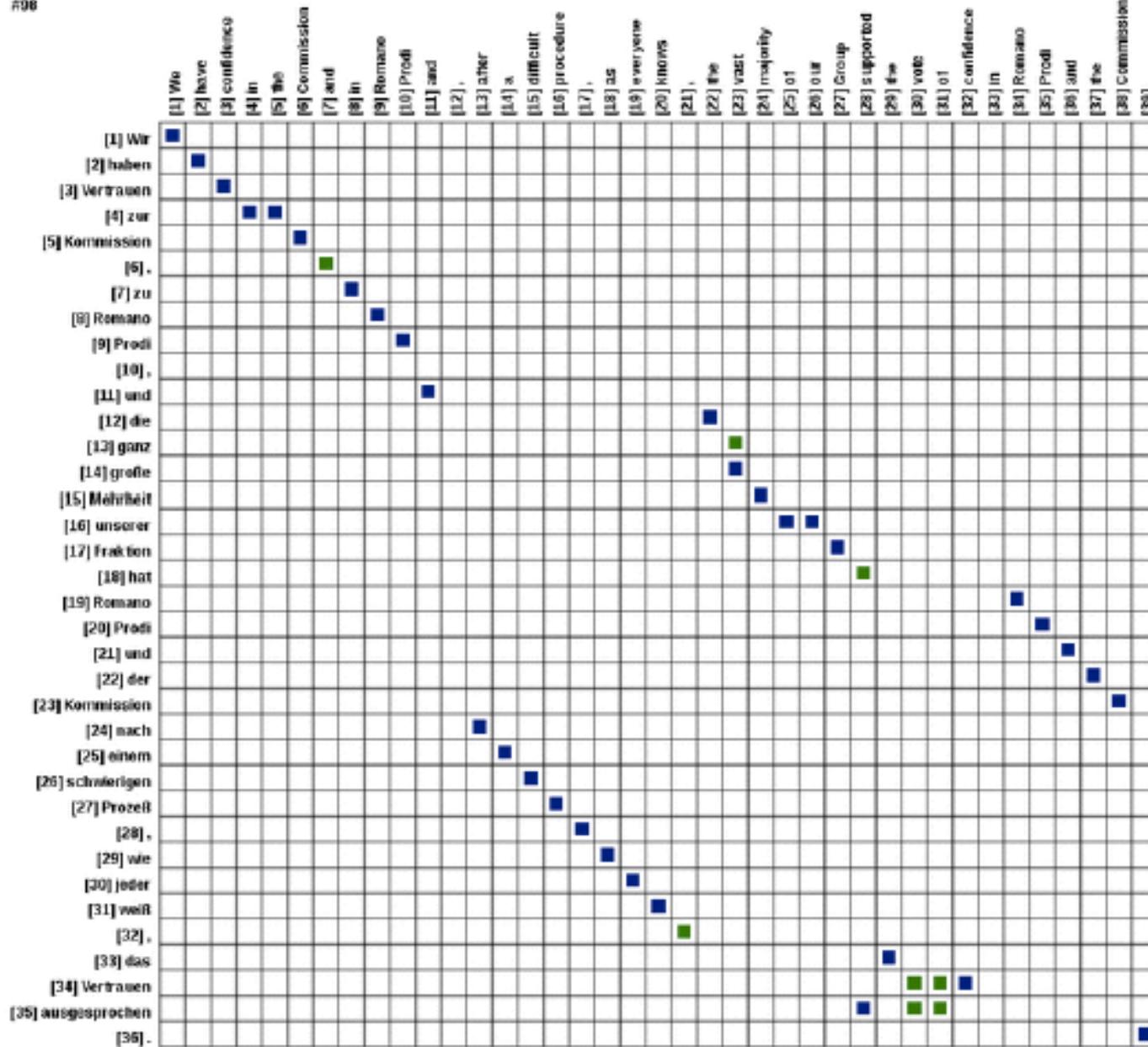
$$\text{Model 2} = \sum_{\mathbf{a} \in [0,n]^m} \prod_{i=1}^m p(a_i | i, m, n) \times p(e_i | f_{a_i})$$

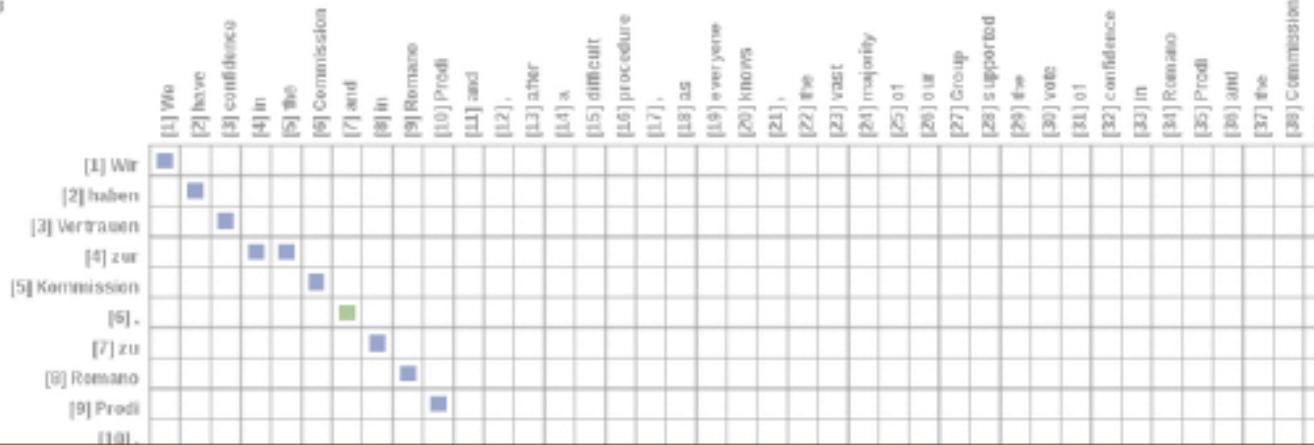
*How much do we know  
when we only know the  
source & target lengths  
and the current position?*











$$p(\mathbf{e} \mid \mathbf{f}, m) = \sum_{\mathbf{a} \in [0, n]^m} \prod_{i=1}^m p(a_i) \times p(e_i \mid f_{a_i})$$

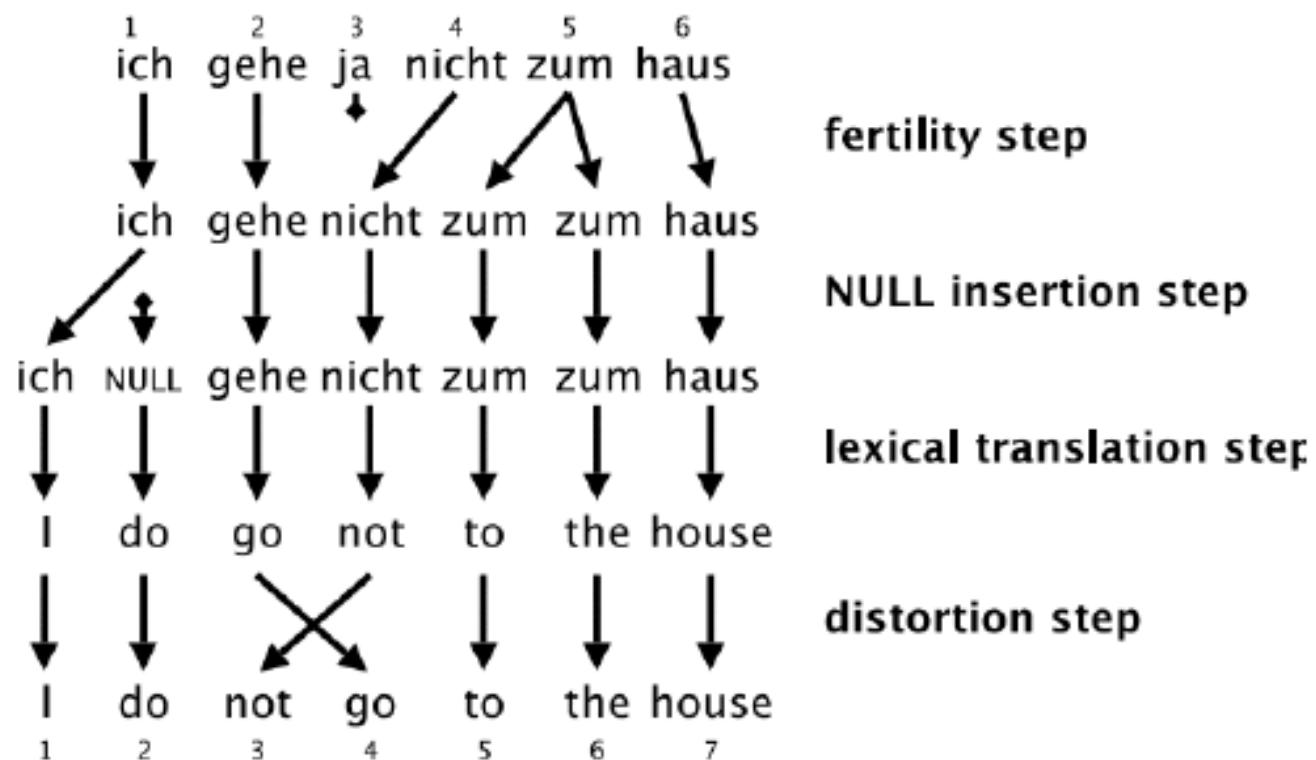
$$\text{Model 2} = \sum_{\mathbf{a} \in [0, n]^m} \prod_{i=1}^m p(a_i \mid i, m, n) \times p(e_i \mid f_{a_i})$$

$$\text{HMM} = \sum_{\mathbf{a} \in [0, n]^m} \prod_{i=1}^m p(a_i \mid a_{i-1}) \times p(e_i \mid f_{a_i})$$

# Fertility Models

- The models we have considered so far have been efficient
- This efficiency has come at a modeling cost:
  - What is to stop the model from “translating” a word 0, 1, 2, or 100 times?
  - We introduce *fertility models* to deal with this

# IBM Model 3

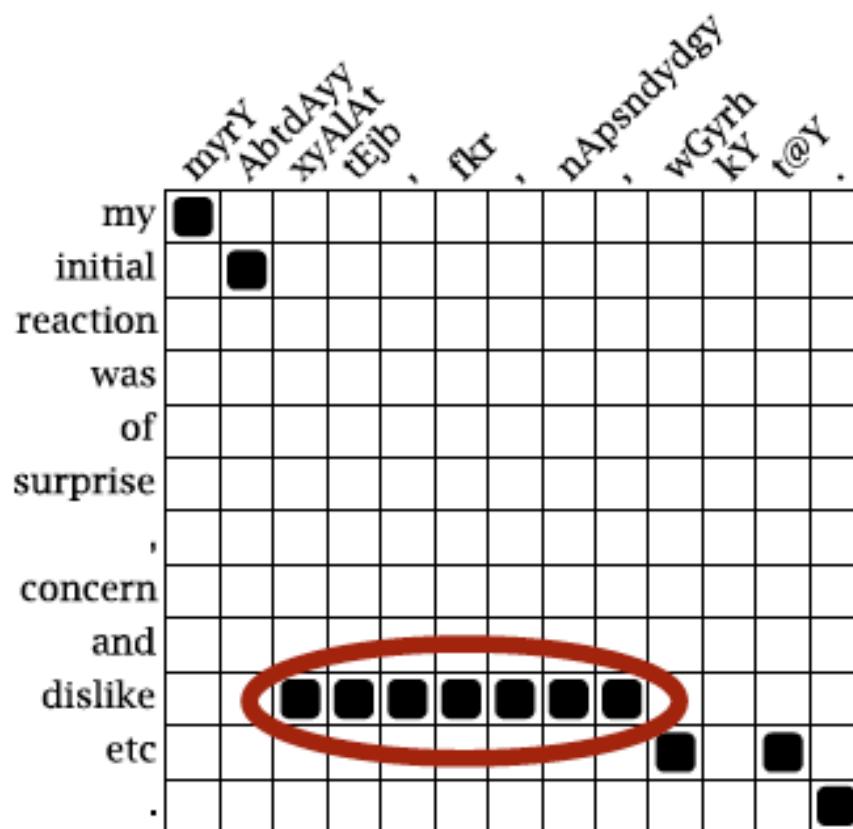


# Fertility

$$p(\mathbf{e} \mid \mathbf{f}, m) = \sum_{\mathbf{a} \in [0,n]^m} p(\mathbf{a} \mid \mathbf{f}, m) \times \prod_{i=1}^m p(e_i \mid f_{a_i})$$

- Fertility models mean that we can no longer exploit conditional independencies to write  $p(\mathbf{a} \mid \mathbf{f}, m)$  as a series of local alignment decisions.
- *How do we compute the statistics required for EM training?*

# Pitfalls of Conditional Models



IBM Model 4 alignment

$p(f|e)$

	michael	geht	davon	aus	dass	er	im	haus	bleibt
michael	█								
assumes		█							
that			█						
he				█					
will					█				
stay						█			
in							█		
the								█	
house									█

English to German

# A few tricks...

# A few tricks...

$p(f|e)$

	michael	geht	davon	aus	dass	er	im	haus	bleibt
michael									
assumes									
that									
he									
will									
stay									
in									
the									
house									

English to German

	michael	geht	davon	aus	dass	er	im	haus	bleibt
michael									
assumes									
that									
he									
will									
stay									
in									
the									
house									

German to English

$p(e|f)$

# A few tricks...

$p(f|e)$

	michael	geht	davon	aus	.	dass	er	im	haus	bleibt
michael										
assumes										
that										
he										
will										
stay										
in										
the										
house										

English to German

	michael	geht	davon	aus	.	dass	er	im	haus	bleibt
michael										
assumes										
that										
he										
will										
stay										
in										
the										
house										

German to English

$p(e|f)$

	michael	geht	davon	aus	.	dass	er	im	haus	bleibt
michael										
assumes										
that										
he										
will										
stay										
in										
the										
house										

Intersection / Union

# Phrase-based Translation Models

# Translational Equivalence

*Er hat die Prüfung **bestanden**, jedoch nur knapp*

He **insisted on** the test, but just barely.

He **passed** the test, but just barely.

How do lexical translation models deal with contextual information?

# Translation model

- With a **latent variable**, we introduce a decomposition into **phrases** which translate **independently**:

$$p(f, a | e) = p(a) \prod_{\langle \bar{e}, \bar{f} \rangle \in a} p(\bar{f} | \bar{e})$$

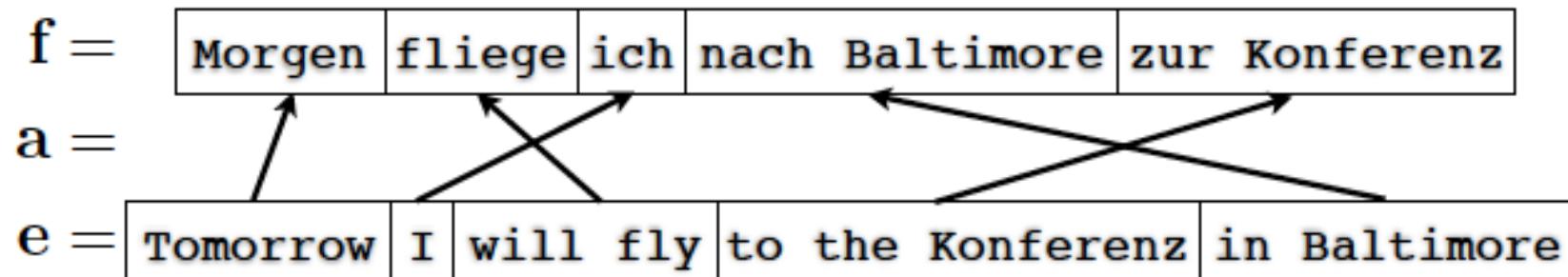
$f =$  Morgen fliege ich nach Baltimore zur Konferenz

$e =$  Tomorrow I will fly to the Konferenz in Baltimore

# Translation model

- With a **latent variable**, we introduce a decomposition into **phrases** which translate **independently**:

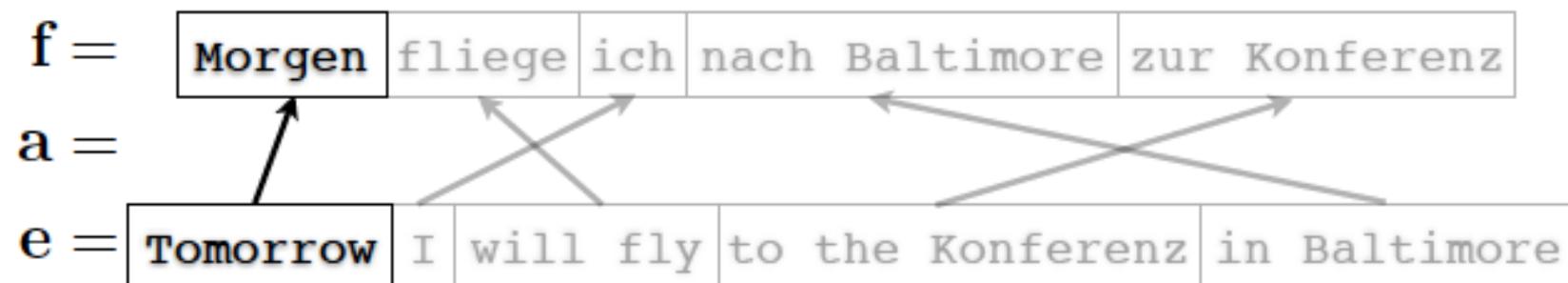
$$p(f, a | e) = p(a) \prod_{\langle \bar{e}, \bar{f} \rangle \in a} p(\bar{f} | \bar{e})$$



# Translation model

- With a **latent variable**, we introduce a decomposition into **phrases** which translate **independently**:

$$p(\mathbf{f}, \mathbf{a} \mid \mathbf{e}) = p(\mathbf{a}) \prod_{\langle \bar{\mathbf{e}}, \bar{\mathbf{f}} \rangle \in \mathbf{a}} p(\bar{\mathbf{f}} \mid \bar{\mathbf{e}})$$

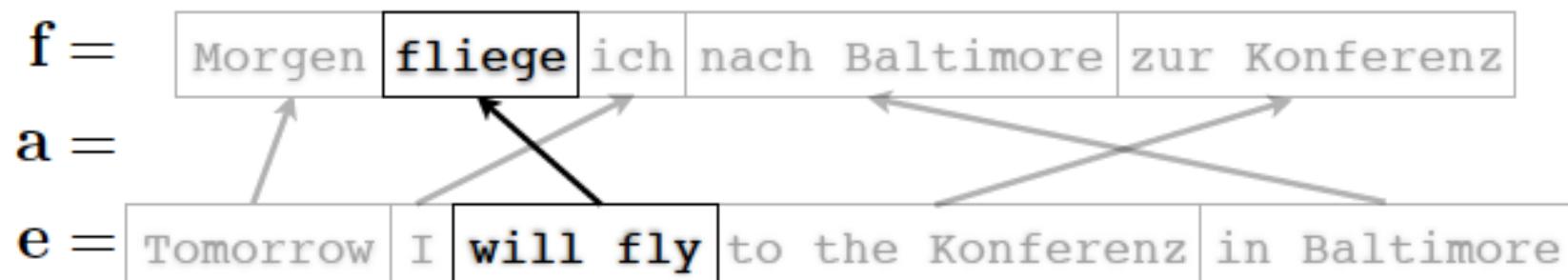


$$p(\text{Morgen} \mid \text{Tomorrow})$$

# Translation model

- With a **latent variable**, we introduce a decomposition into **phrases** which translate **independently**:

$$p(\mathbf{f}, \mathbf{a} \mid \mathbf{e}) = p(\mathbf{a}) \prod_{\langle \bar{\mathbf{e}}, \bar{\mathbf{f}} \rangle \in \mathbf{a}} p(\bar{\mathbf{f}} \mid \bar{\mathbf{e}})$$

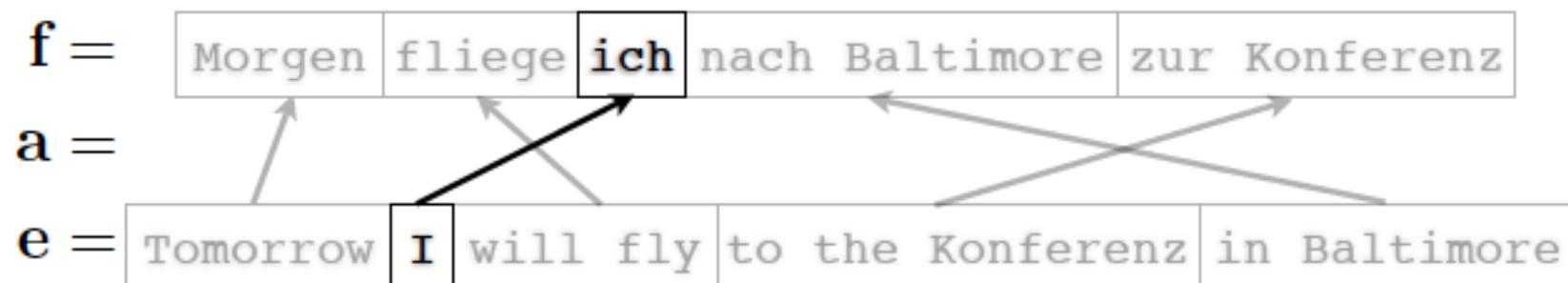


$$p(\text{Morgen} \mid \text{Tomorrow}) \times p(\text{fliege} \mid \text{will fly})$$

# Translation model

- With a **latent variable**, we introduce a decomposition into **phrases** which translate **independently**:

$$p(\mathbf{f}, \mathbf{a} \mid \mathbf{e}) = p(\mathbf{a}) \prod_{\langle \bar{\mathbf{e}}, \bar{\mathbf{f}} \rangle \in \mathbf{a}} p(\bar{\mathbf{f}} \mid \bar{\mathbf{e}})$$

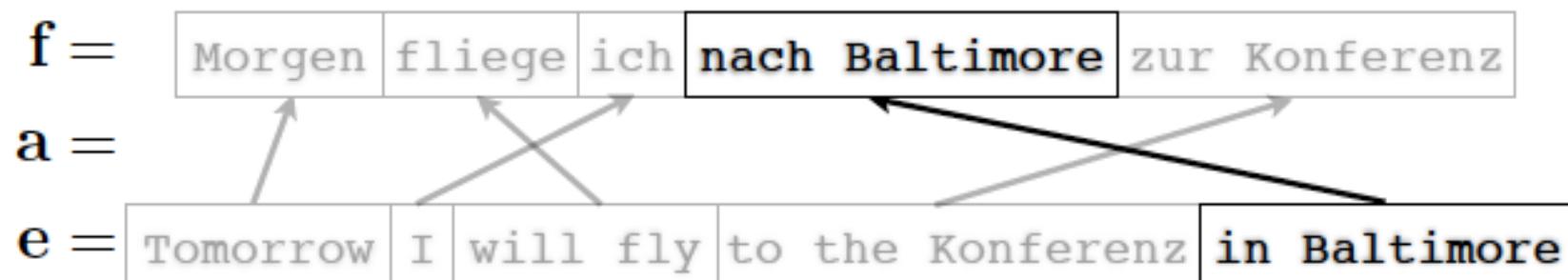


$$p(\text{Morgen} \mid \text{Tomorrow}) \times p(\text{fliege} \mid \text{will fly}) \times p(\text{ich} \mid \text{I})$$

# Translation model

- With a **latent variable**, we introduce a decomposition into **phrases** which translate **independently**:

$$p(\mathbf{f}, \mathbf{a} \mid \mathbf{e}) = p(\mathbf{a}) \prod_{\langle \bar{\mathbf{e}}, \bar{\mathbf{f}} \rangle \in \mathbf{a}} p(\bar{\mathbf{f}} \mid \bar{\mathbf{e}})$$



$$p(\text{Morgen} \mid \text{Tomorrow}) \times p(\text{fliege} \mid \text{will fly}) \times p(\text{ich} \mid \text{I}) \times \dots$$

# Translation model

- With a **latent variable**, we introduce a decomposition into **phrases** which translate **independently**:

$$p(\mathbf{f}, \mathbf{a} \mid \mathbf{e}) = p(\mathbf{a}) \prod_{\langle \bar{\mathbf{e}}, \bar{\mathbf{f}} \rangle \in \mathbf{a}} p(\bar{\mathbf{f}} \mid \bar{\mathbf{e}})$$

Marginalize to get  $p(\mathbf{f} \mid \mathbf{e})$ :

$$p(\mathbf{f} \mid \mathbf{e}) = \sum_{\mathbf{a} \in \mathcal{A}} p(\mathbf{a}) \prod_{\langle \bar{\mathbf{e}}, \bar{\mathbf{f}} \rangle \in \mathbf{a}} p(\bar{\mathbf{f}} \mid \bar{\mathbf{e}})$$

# Phrases

- Contiguous strings of words
- Phrases are not necessarily syntactic constituents
- Usually have maximum limits
- Phrases subsume words (words are phrases)

# Phrase Tables

	$\bar{f}$	$\bar{e}$	$p(\bar{f}   \bar{e})$
das Thema		the issue	0.41
		the point	0.72
		the subject	0.47
		the thema	0.99
es gibt		there is	0.96
		there are	0.72
morgen		tomorrow	0.9
fliege ich		will I fly	0.63
		will fly	0.17
		I will fly	0.13

$$P(a)$$

- Two responsibilities
  - Divide the source sentence into phrases
    - Standard approach: uniform distribution over all possible segmentations
    - How many segmentations are there?
  - Reorder the phrases
    - Standard approach: Markov model on phrases (parameterized with log-linear model)

# Learning Phrases

- Latent segmentation variable
- Latent phrasal inventory
- Parallel data
- EM?

Computational problem: summing over all segmentations and alignments is #P-complete

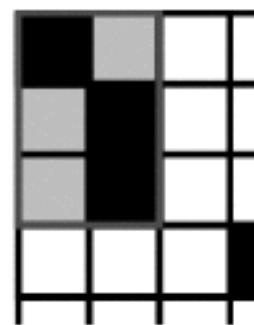
$$p(f | e) = \sum_{a \in \mathcal{A}} p(a) \prod_{(\bar{e}, \bar{f}) \in a} p(\bar{f} | \bar{e})$$

Modeling problem: MLE has a degenerate solution.

# Learning Phrases

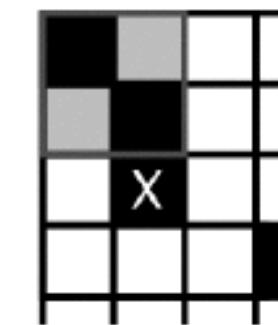
- Three stages
  - word alignment
  - extraction of phrases
  - estimation of phrase probabilities

# Consistent Phrases



consistent

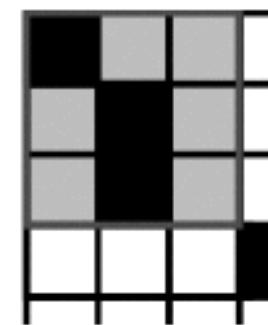
**ok**



inconsistent

**violated**

one alignment  
point outside



consistent

**ok**

unaligned  
word is fine

All words of the phrase pair have to align to each other.

# Scoring Phrase Pairs

- Phrase pair extraction: collect all phrase pairs from the data
- Phrase pair scoring: assign probabilities to phrase translations
- Score by relative frequency:

$$\phi(\bar{f}|\bar{e}) = \frac{\text{count}(\bar{e}, \bar{f})}{\sum_{\bar{f}_i} \text{count}(\bar{e}, \bar{f}_i)}$$

# Decoding in Phrase-based Models

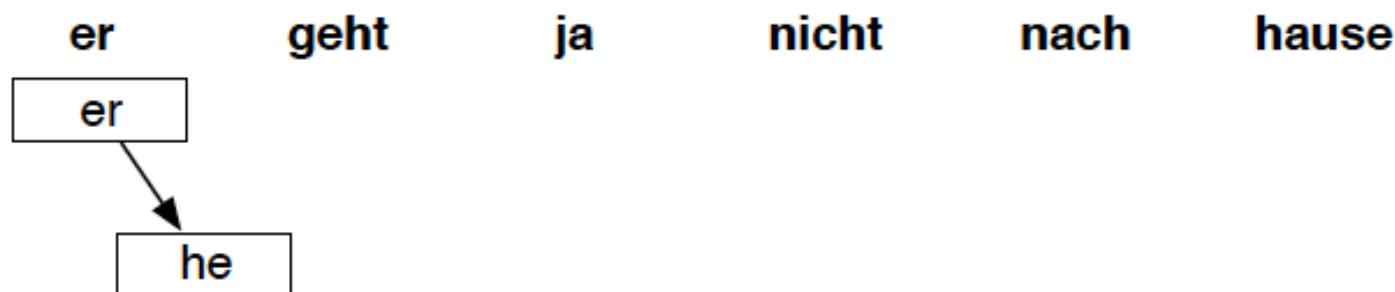
## **Translation Process**

- Task: translate this sentence from German into English

er            geht            ja            nicht            nach            hause

## Translation Process

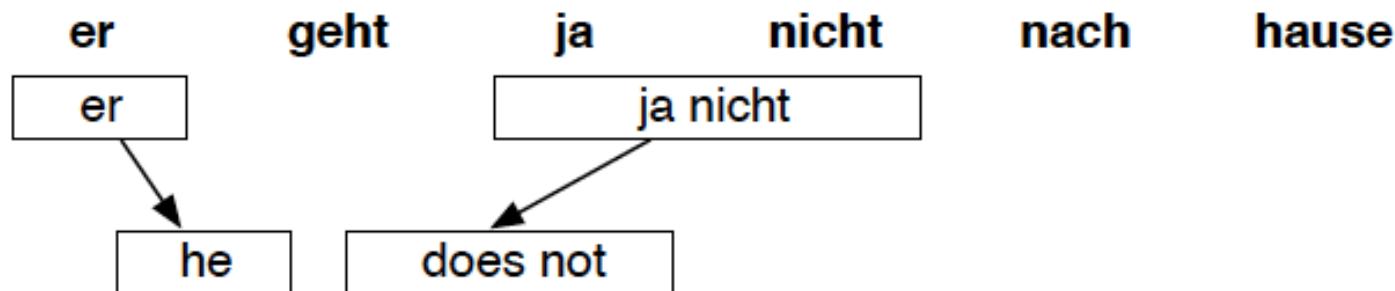
- Task: translate this sentence from German into English



- Pick phrase in input, translate

## Translation Process

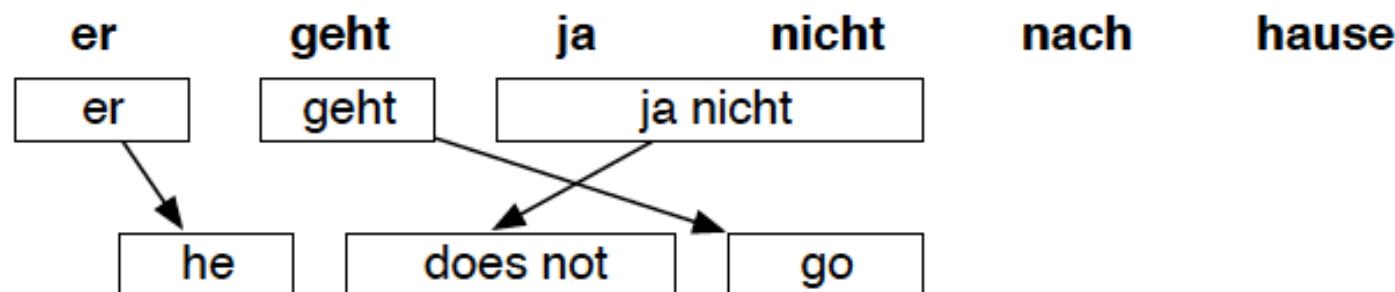
- Task: translate this sentence from German into English



- Pick phrase in input, translate
  - it is allowed to pick words out of sequence reordering
  - phrases may have multiple words: many-to-many translation

## Translation Process

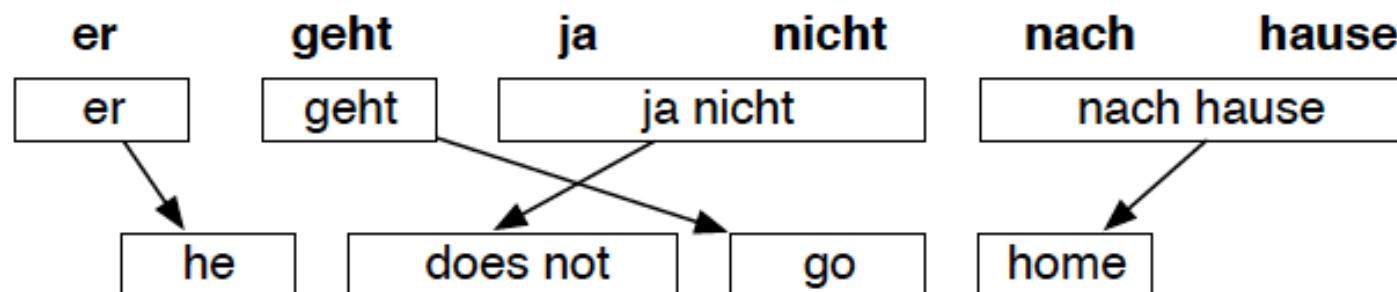
- Task: translate this sentence from German into English



- Pick phrase in input, translate

## Translation Process

- Task: translate this sentence from German into English



- Pick phrase in input, translate

# Computing Translation Probability

- Probabilistic model for phrase-based translation:

$$e_{\text{best}} = \operatorname{argmax}_e \prod_{i=1}^I \phi(\bar{f}_i | \bar{e}_i) d(start_i - end_{i-1} - 1) p_{\text{LM}}(e)$$

- Score is computed incrementally for each partial hypothesis
- Components

**Phrase translation** Picking phrase  $\bar{f}_i$  to be translated as a phrase  $\bar{e}_i$

→ look up score  $\phi(\bar{f}_i | \bar{e}_i)$  from phrase translation table

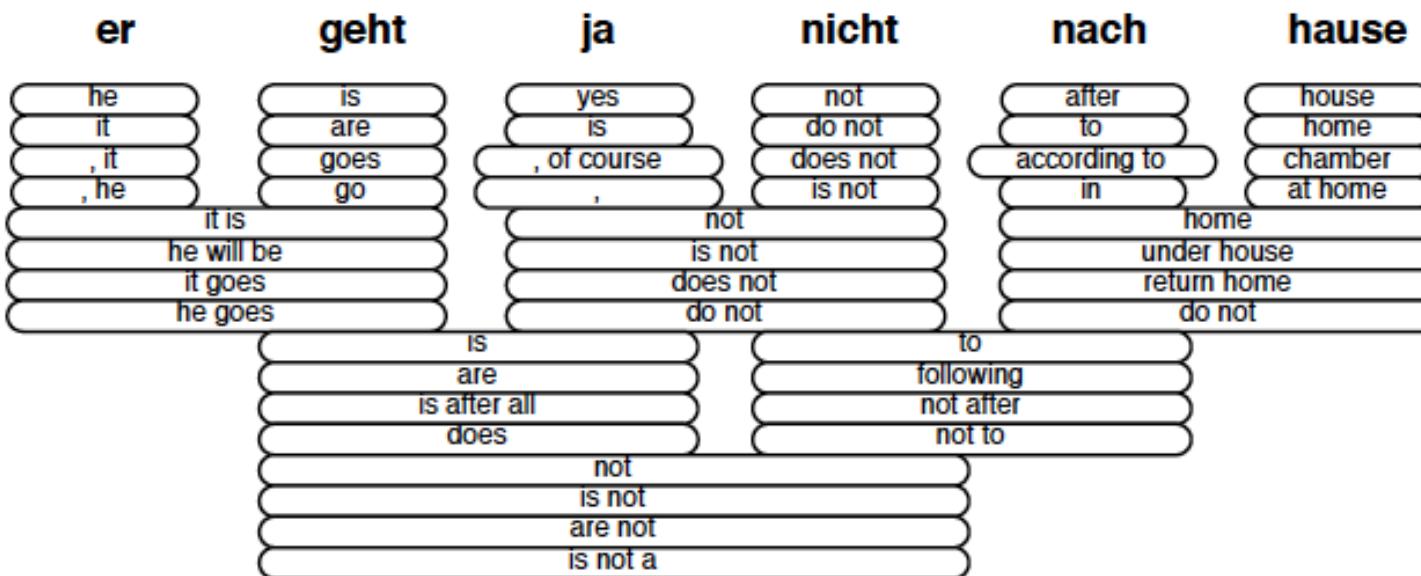
**Reordering** Previous phrase ended in  $end_{i-1}$ , current phrase starts at  $start_i$

→ compute  $d(start_i - end_{i-1} - 1)$

**Language model** For  $n$ -gram model, need to keep track of last  $n - 1$  words

→ compute score  $p_{\text{LM}}(w_i | w_{i-(n-1)}, \dots, w_{i-1})$  for added words  $w_i$

## Translation Options



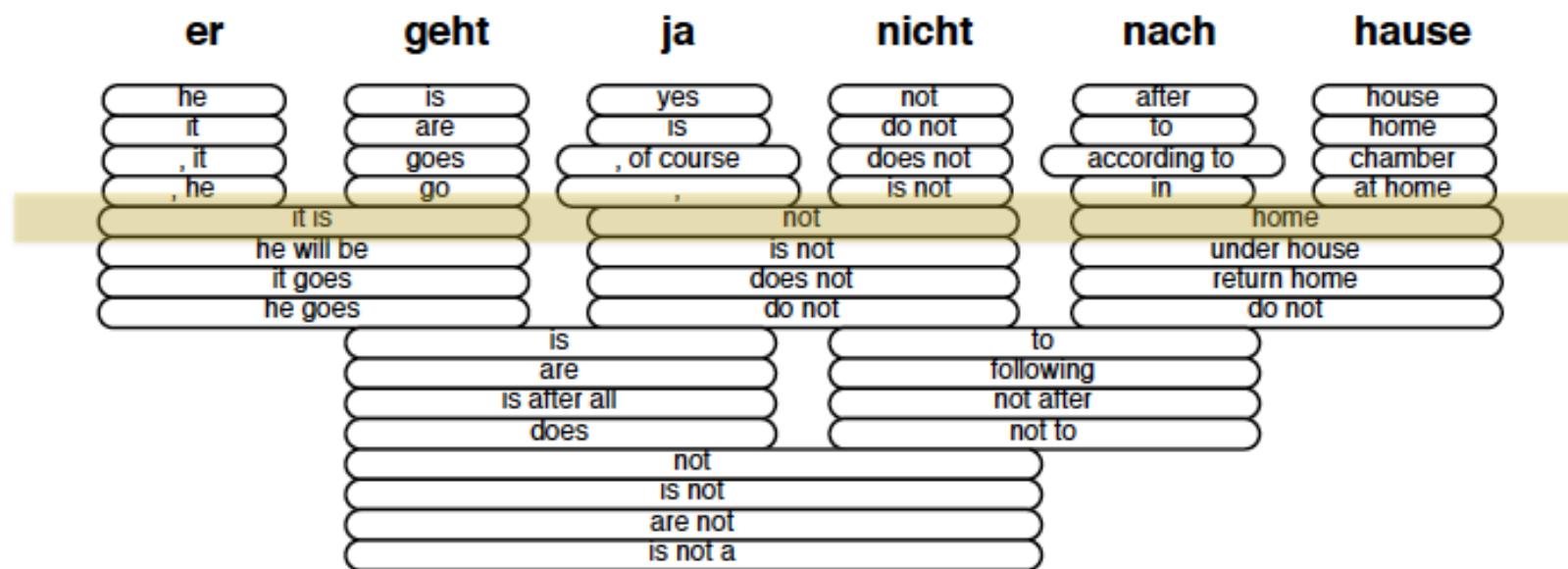
- Many translation options to choose from
  - in Europarl phrase table: 2727 matching phrase pairs for this sentence
  - by pruning to the top 20 per phrase, 202 translation options remain

## Translation Options

er	geht	ja	nicht	nach	hause
he	is	yes	not	after	house
it	are	is	do not	to	home
, it	goes	, of course	does not	according to	chamber
, he	go	,	is not	in	at home
it is		not		home	
he will be		is not		under house	
it goes		does not		return home	
he goes		do not		do not	
	IS		to		
	are		following		
	is after all		not after		
	does		not to		
	not				
	is not				
	are not				
	is not a				

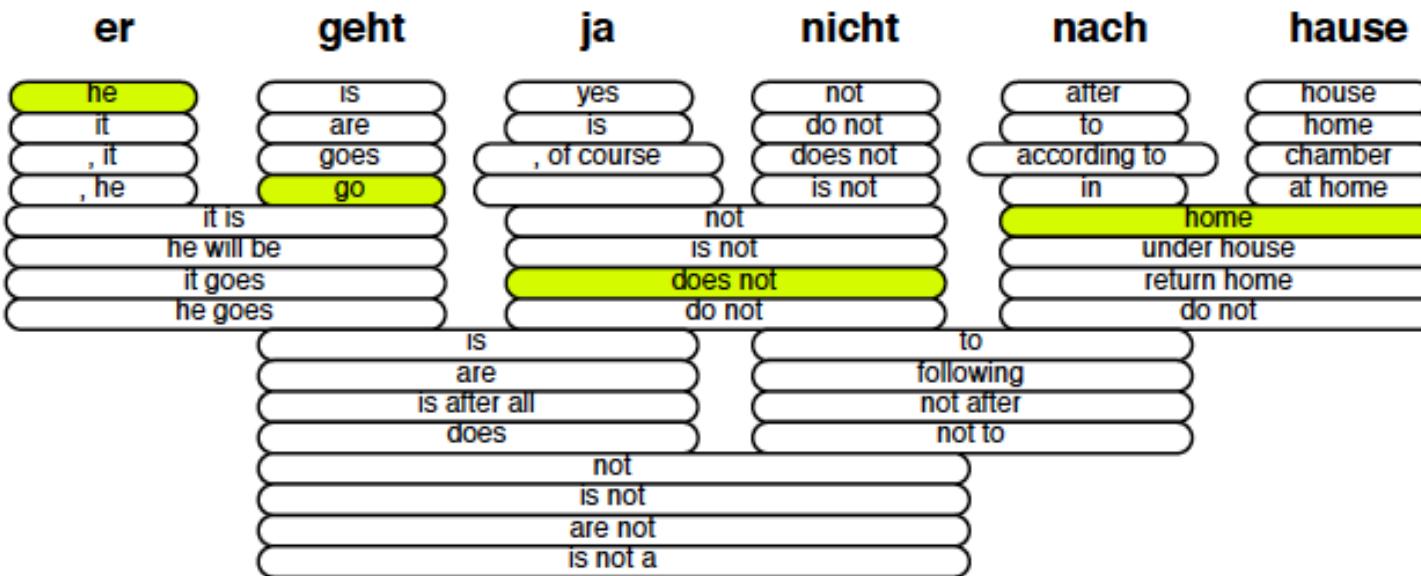
- Many translation options to choose from
  - in Europarl phrase table: 2727 matching phrase pairs for this sentence
  - by pruning to the top 20 per phrase, 202 translation options remain

## Translation Options



- Many translation options to choose from
  - in Europarl phrase table: 2727 matching phrase pairs for this sentence
  - by pruning to the top 20 per phrase, 202 translation options remain

## Translation Options



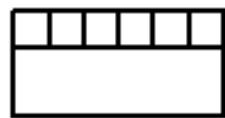
- The machine translation decoder does not know the right answer
    - picking the right translation options
    - arranging them in the right order
- Search problem solved by heuristic beam search

# Decoding: Precompute Translation Options



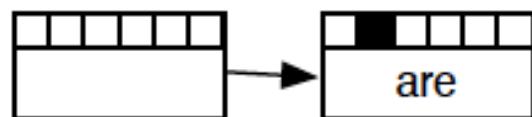
consult phrase translation table for all input phrases

# Decoding: Start with Initial Hypothesis



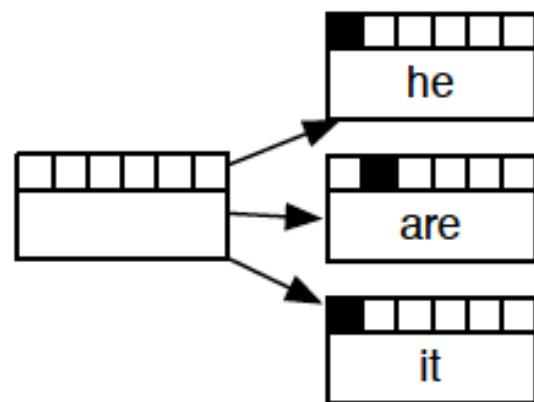
initial hypothesis: no input words covered, no output produced

# Decoding: Hypothesis Expansion



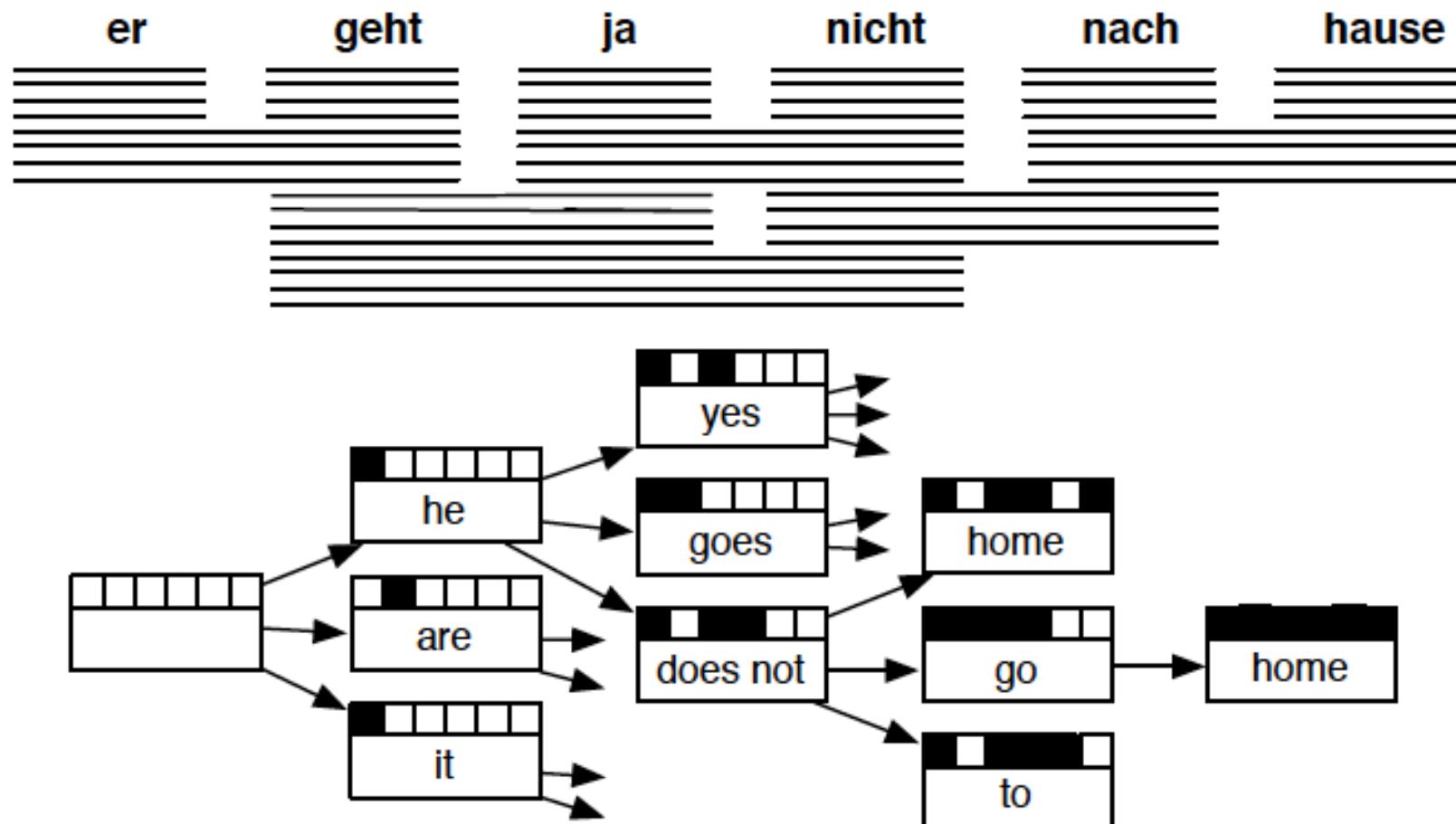
pick any translation option, create new hypothesis

# Decoding: Hypothesis Expansion



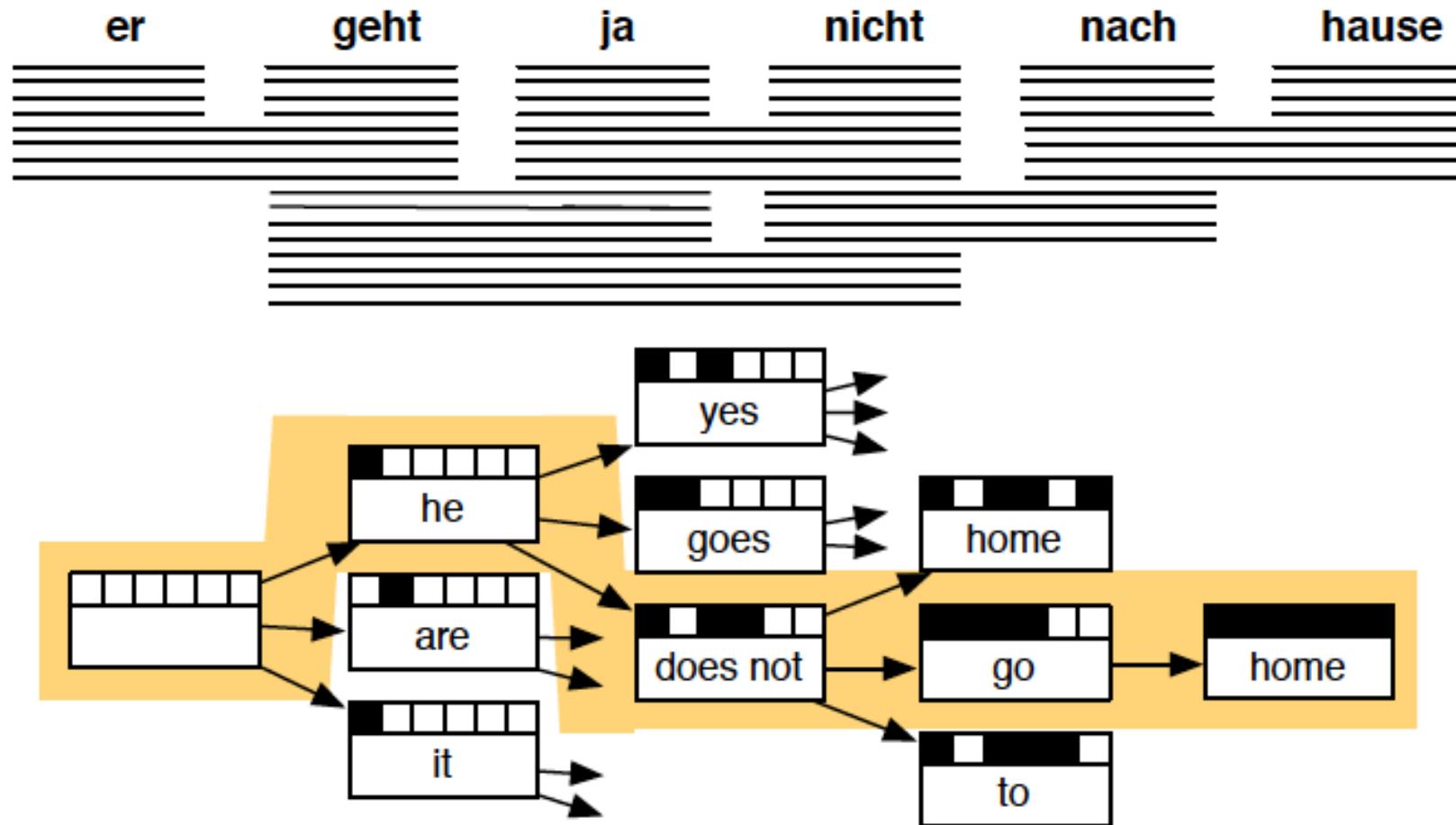
create hypotheses for all other translation options

# Decoding: Hypothesis Expansion



also create hypotheses from created partial hypothesis

# Decoding: Find Best Path



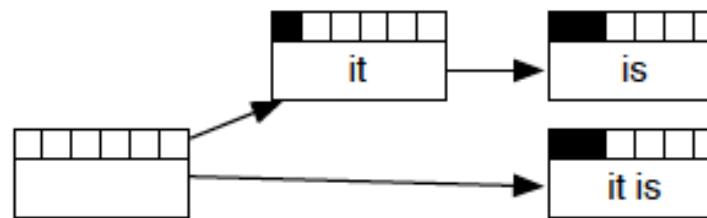
backtrack from highest scoring complete hypothesis

# Size of the search space

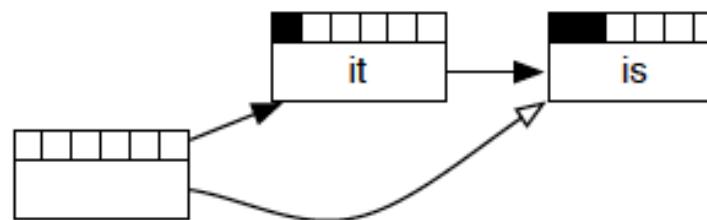
- Decoding is an NP-complete problem
- The search space is exponential
- How can we reduce the search?
- Dynamic programming (risk free)
- Heuristic search (risky)

# Recombination

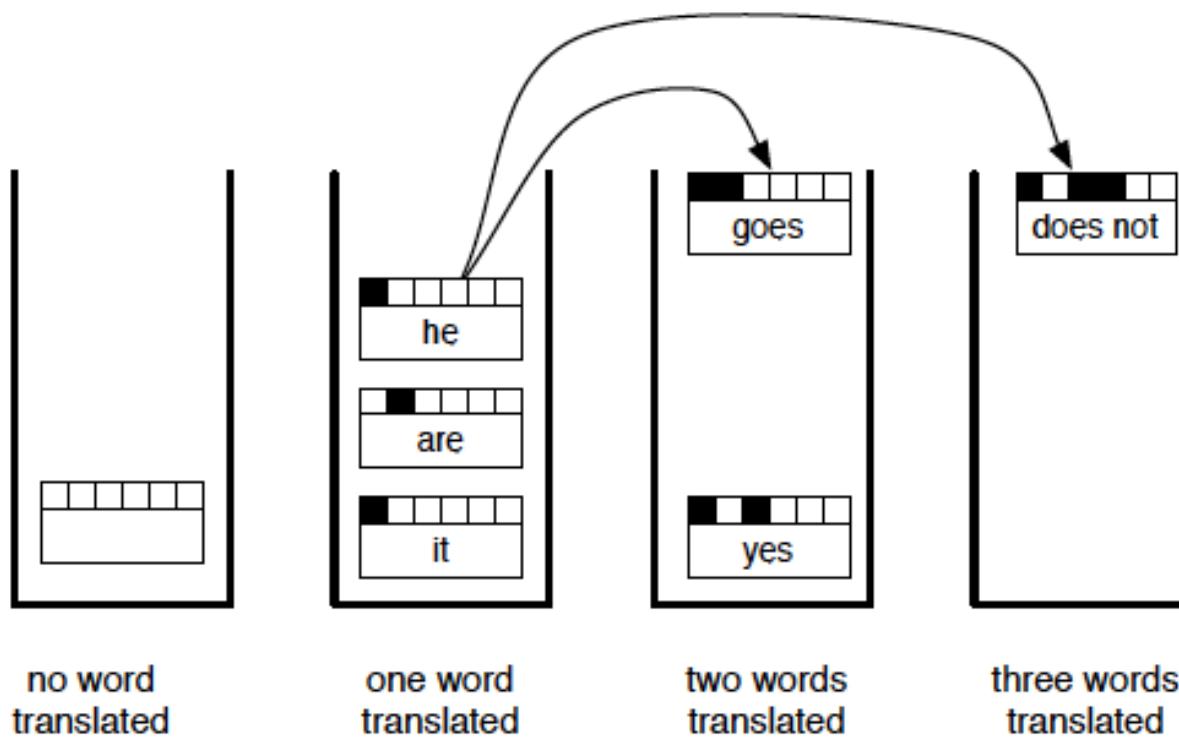
- Two hypothesis paths lead to two matching hypotheses
  - same number of foreign words translated
  - same English words in the output
  - different scores



- Worse hypothesis is dropped



# Stacks



- Hypothesis expansion in a stack decoder
  - translation option is applied to hypothesis
  - new hypothesis is dropped into a stack further down

# Search Errors

- We are using a **heuristic search** to prune the search space
- There are no guarantees of admissibility (like in A\* search)
- We may therefore prune out a partial hypothesis that would have lead to the most probable translation, if we hadn't pruned it early on

# Model Errors

- Model errors are different than search errors
- If the model scores an incorrect translation higher than higher than the correct translation, then it is the model's fault not the search strategy's fault

# Other Decoding Algorithms

- A\* search
- Greedy hill-climbing
- Using finite state transducers (standard toolkits)

# Discriminative Modeling

# Noisy Channels Again

$$\begin{aligned} \mathbf{e}^* &= \arg \max_{\mathbf{e}} p(\mathbf{e} \mid \mathbf{g}) \\ &= \arg \max_{\mathbf{e}} \frac{p(\mathbf{g} \mid \mathbf{e}) \times p(\mathbf{e})}{p(\mathbf{g})} \\ &= \arg \max_{\mathbf{e}} p(\mathbf{g} \mid \mathbf{e}) \times p(\mathbf{e}) \\ &= \arg \max_{\mathbf{e}} \log p(\mathbf{g} \mid \mathbf{e}) + \log p(\mathbf{e}) \end{aligned}$$

# Noisy Channels Again

$$e^* = \arg \max_e p(e | g)$$

$$= \arg \max_e \frac{p(g | e) \times p(e)}{p(g)}$$

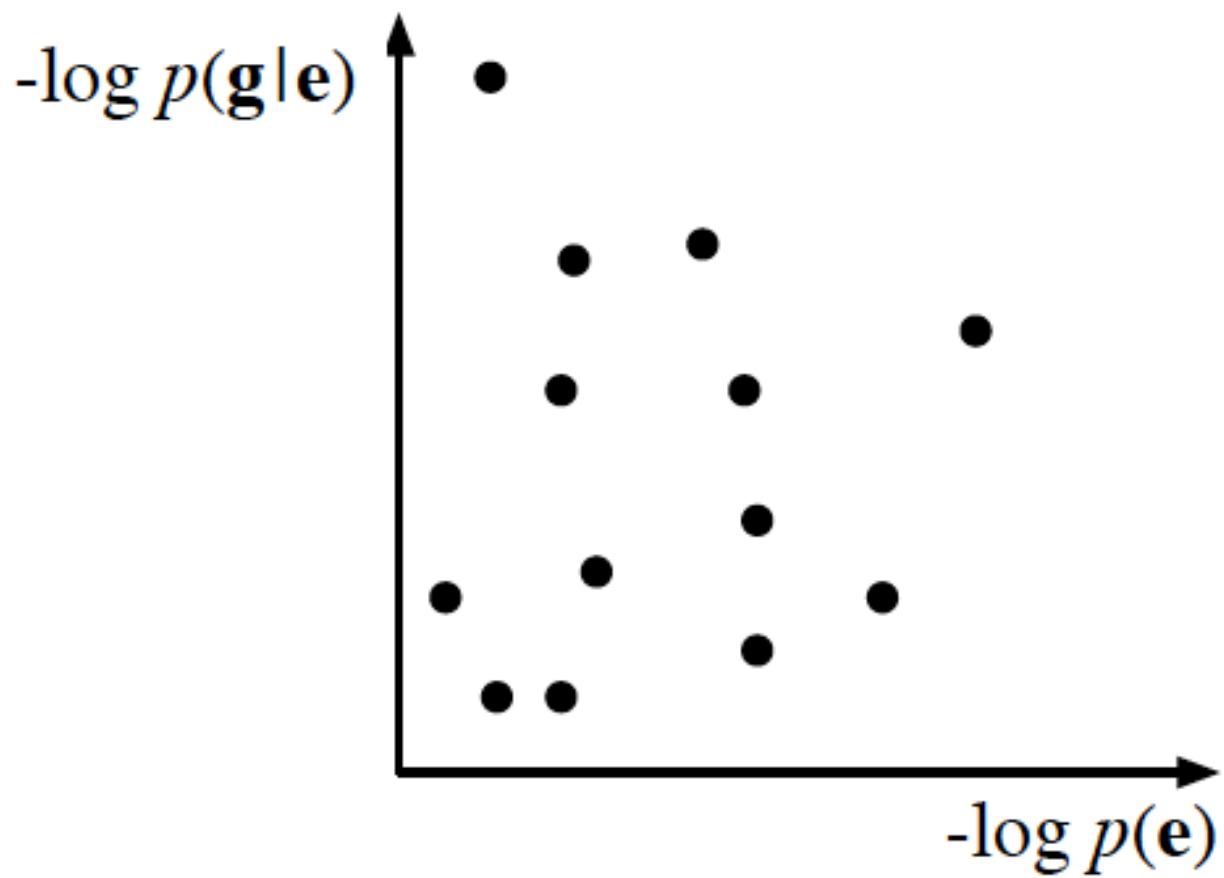
$$= \arg \max_e p(g | e) \times p(e)$$

$$= \arg \max_e \log p(g | e) + \log p(e)$$

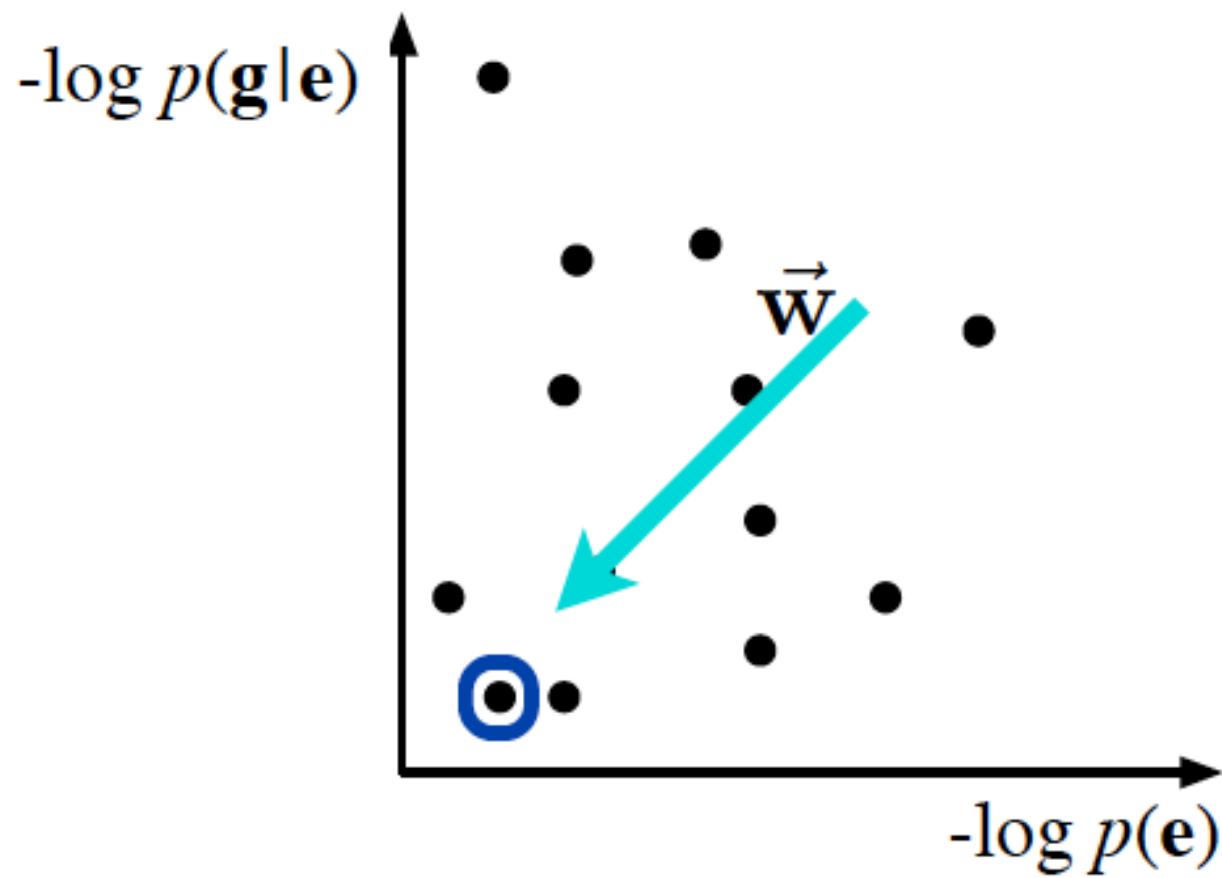
$$= \arg \max_e \underbrace{\begin{bmatrix} 1 \\ 1 \end{bmatrix}}_{w^\top}^\top \underbrace{\begin{bmatrix} \log p(g | e) \\ \log p(e) \end{bmatrix}}_{h(g, e)}$$

**Log-linear Model**

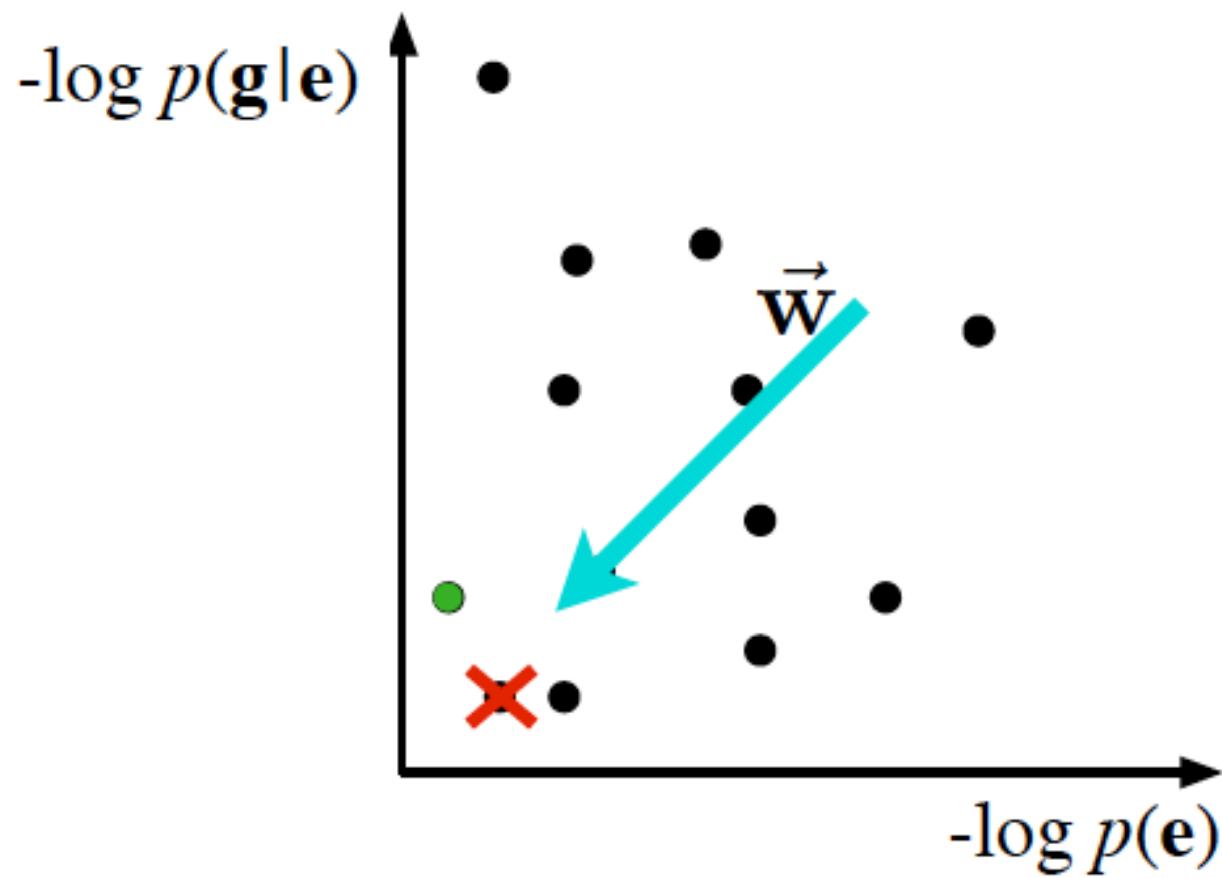
# The Noisy Channel



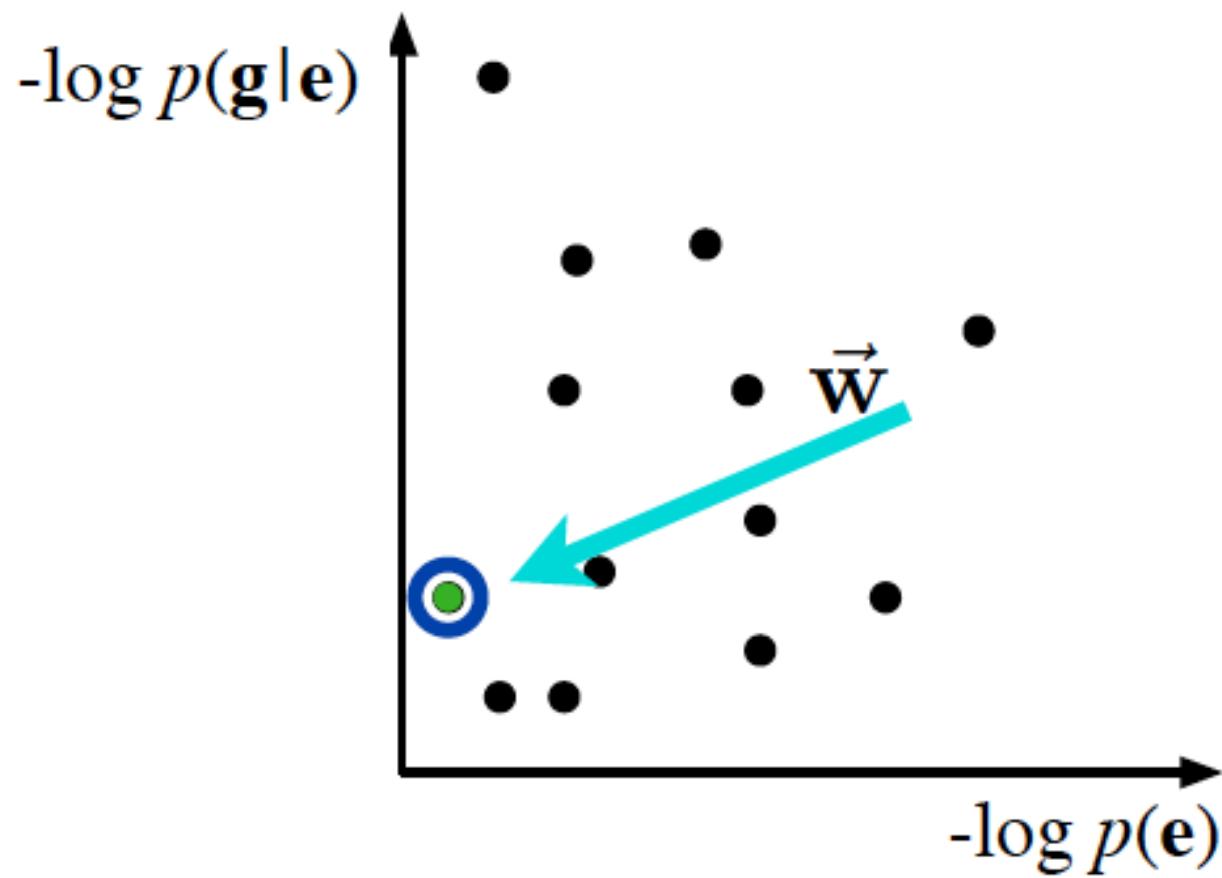
# As a Linear Model



# As a Linear Model



# As a Linear Model

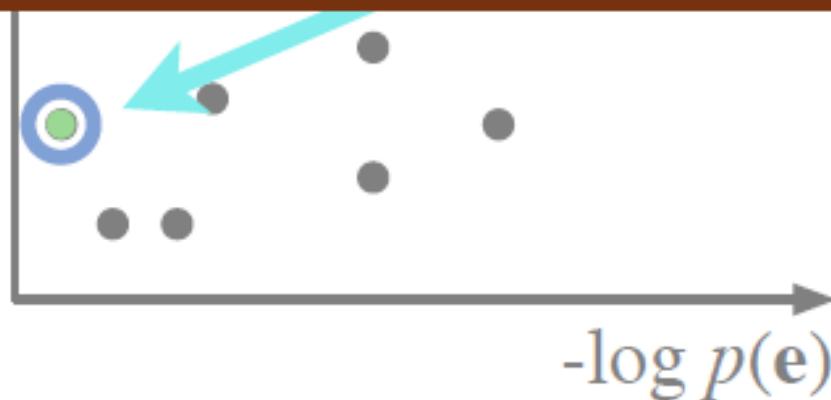


# As a Linear Model

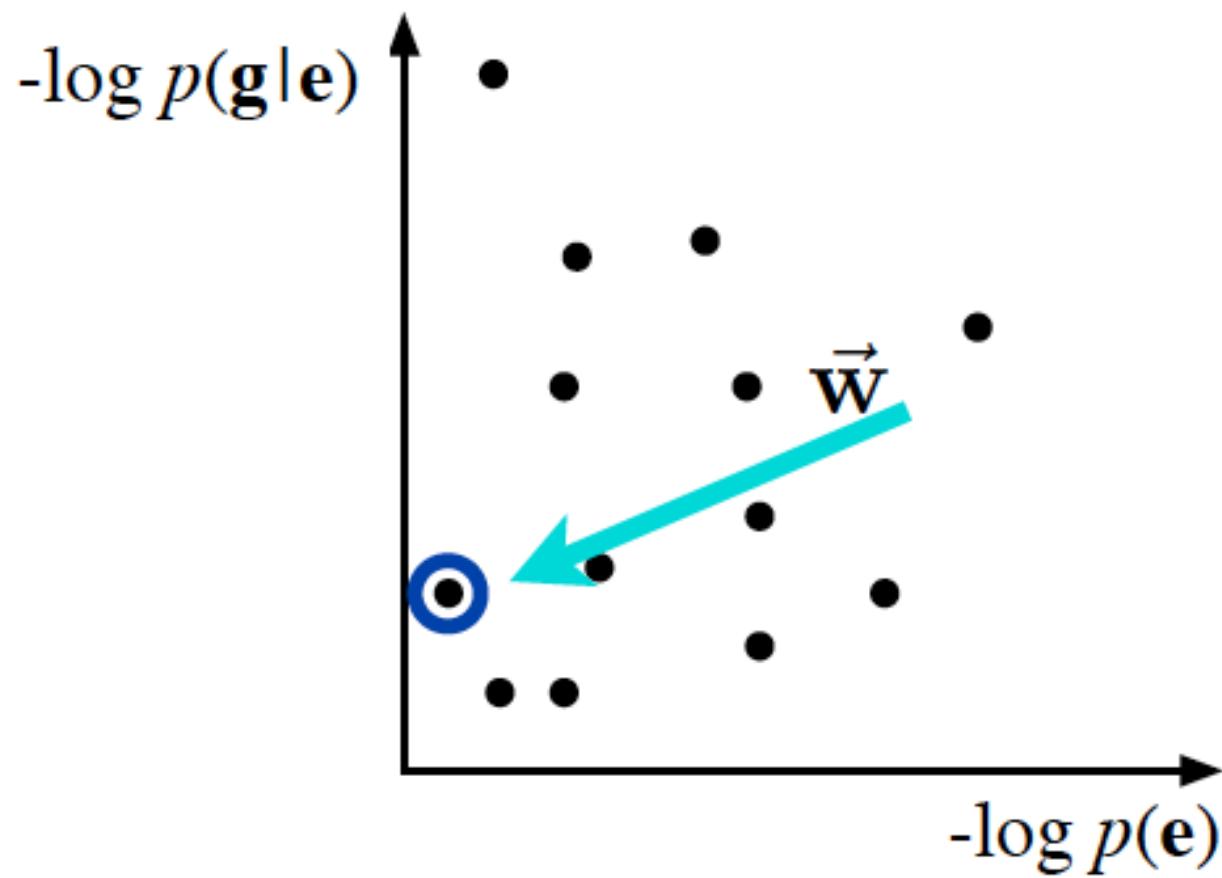
$-\log p(\mathbf{g}|\mathbf{e}) \uparrow$

Improvement I:

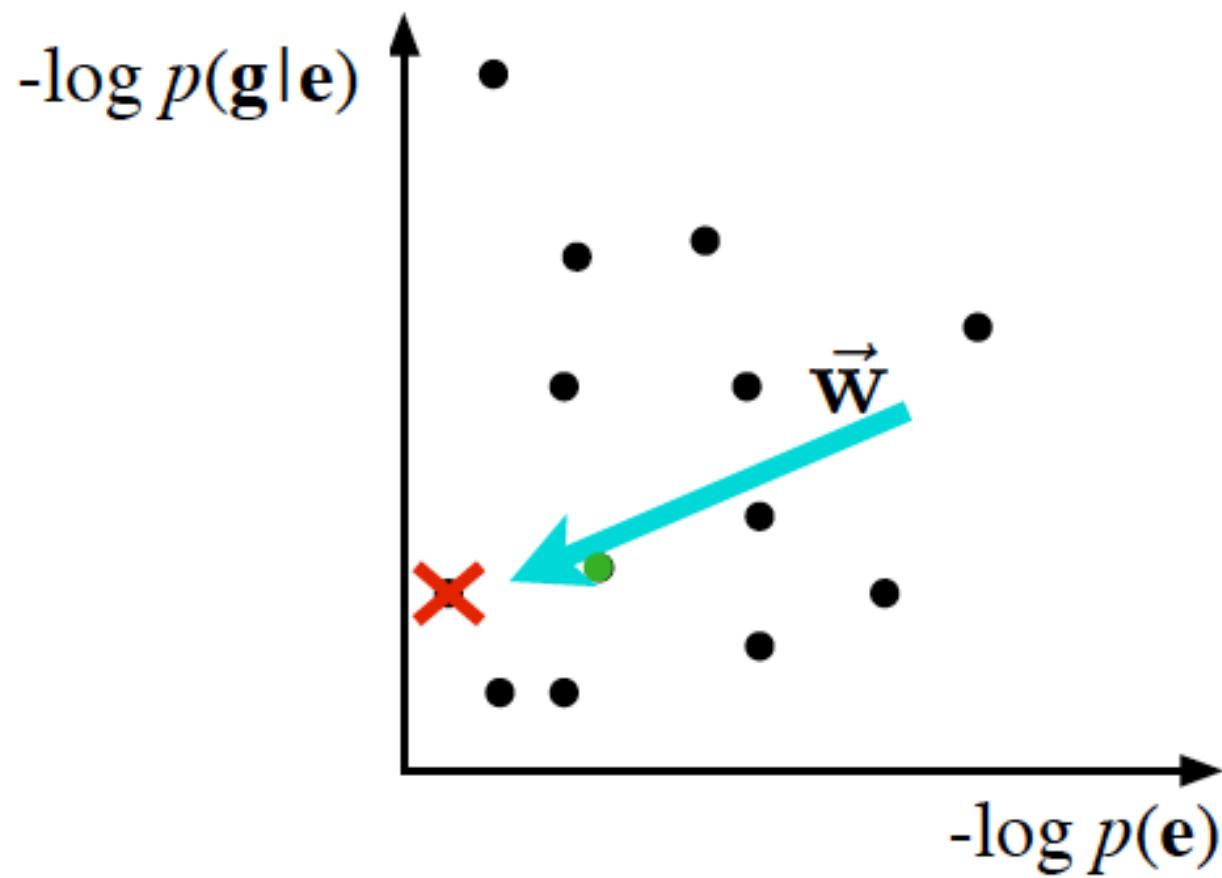
change  $\vec{w}$  to find better translations



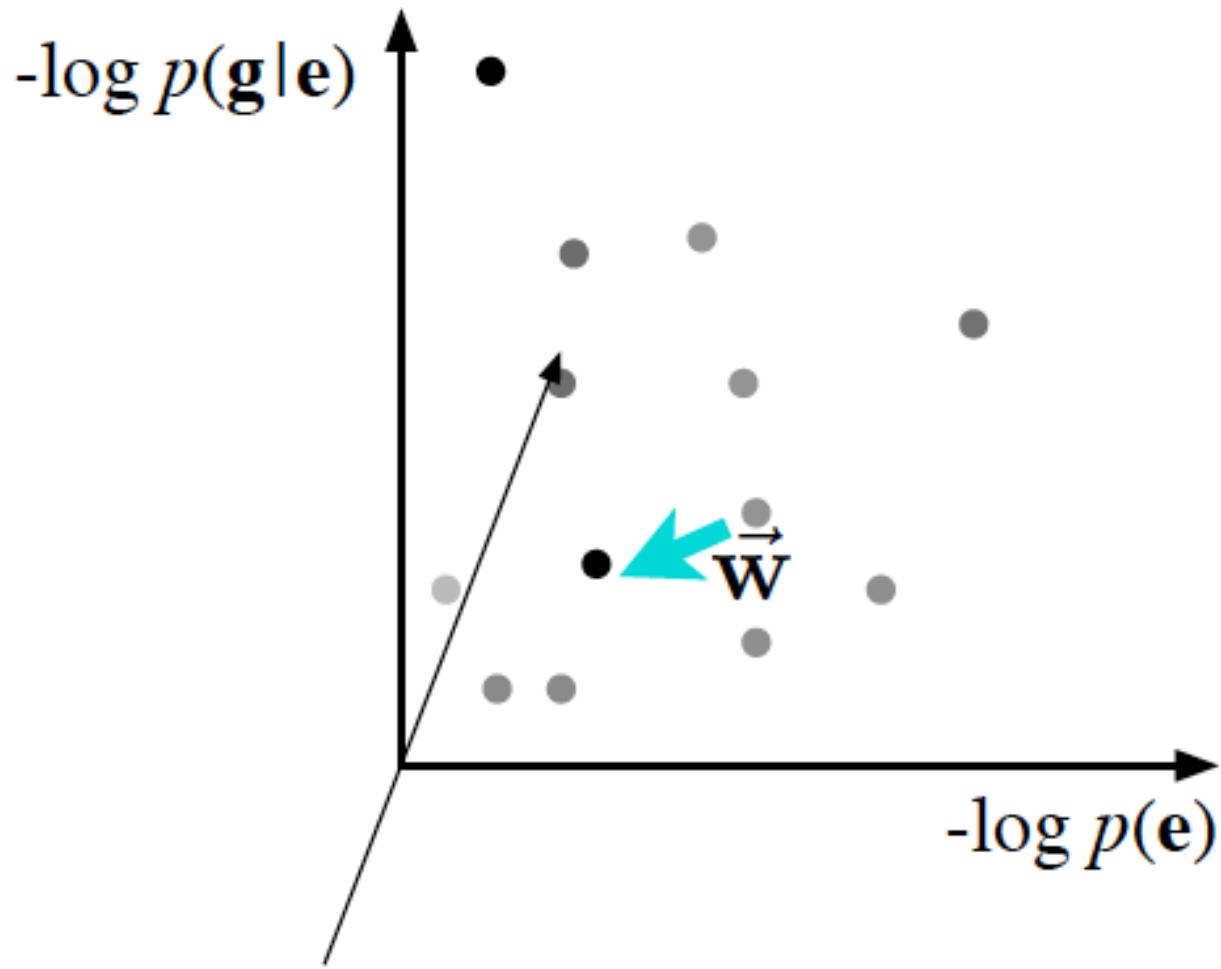
# As a Linear Model



# As a Linear Model



# As a Linear Model



# As a Linear Model

$-\log p(\mathbf{g}|\mathbf{e})$

Improvement 2:

Add dimensions to make points separable



# Linear Models

$$e^* = \arg \max_e w^\top h(g, e)$$

- Improve the modeling capacity of the noisy channel in two ways
  - Reorient the weight vector
  - Add new dimensions (*new features*)
- Questions
  - What features?  $h(g, e)$
  - How do we set the weights?  $w$

# Standard Features

- Target side features
  - $\log p(e)$  [ *n*-gram language model ]
  - Number of words in hypothesis
  - Non-English character count
- Source + target features
  - log relative frequency  $e|f$  of each rule [  $\log \#(e,f) - \log \#(f)$  ]
  - log relative frequency  $f|e$  of each rule [  $\log \#(e,f) - \log \#(e)$  ]
  - “lexical translation” log probability  $e|f$  of each rule [  $\approx \log P_{\text{modelI}}(e|f)$  ]
  - “lexical translation” log probability  $f|e$  of each rule [  $\approx \log P_{\text{modelI}}(f|e)$  ]
- Other features
  - Count of rules/phrases used
  - Reordering pattern probabilities

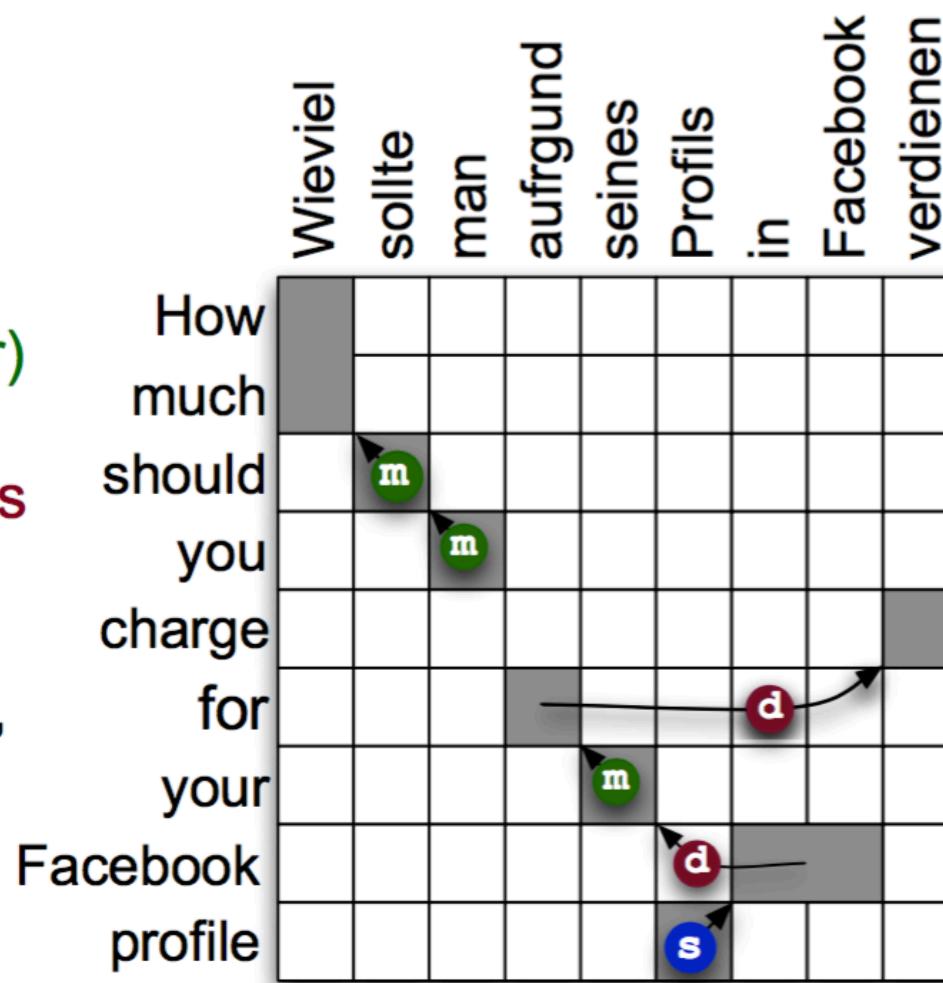
# Lexicalized Reordering

**m:** monotone (keep order)

**s:** swap order

**d:** become discontinuous

Reordering features are probability estimates of s, d, and m



# Learning Weights

- Try to match the reference translation *exactly*
  - Conditional random field
    - Maximize the conditional probability of the reference translations
    - “Average” over the different latent variables
- These methods give “full credit” when the model *exactly* produces the reference and no credit otherwise
- **What is the problem with this?**

# Cost-Sensitive Training

- Assume we have a **cost function** that gives a score for how good/bad a translation is

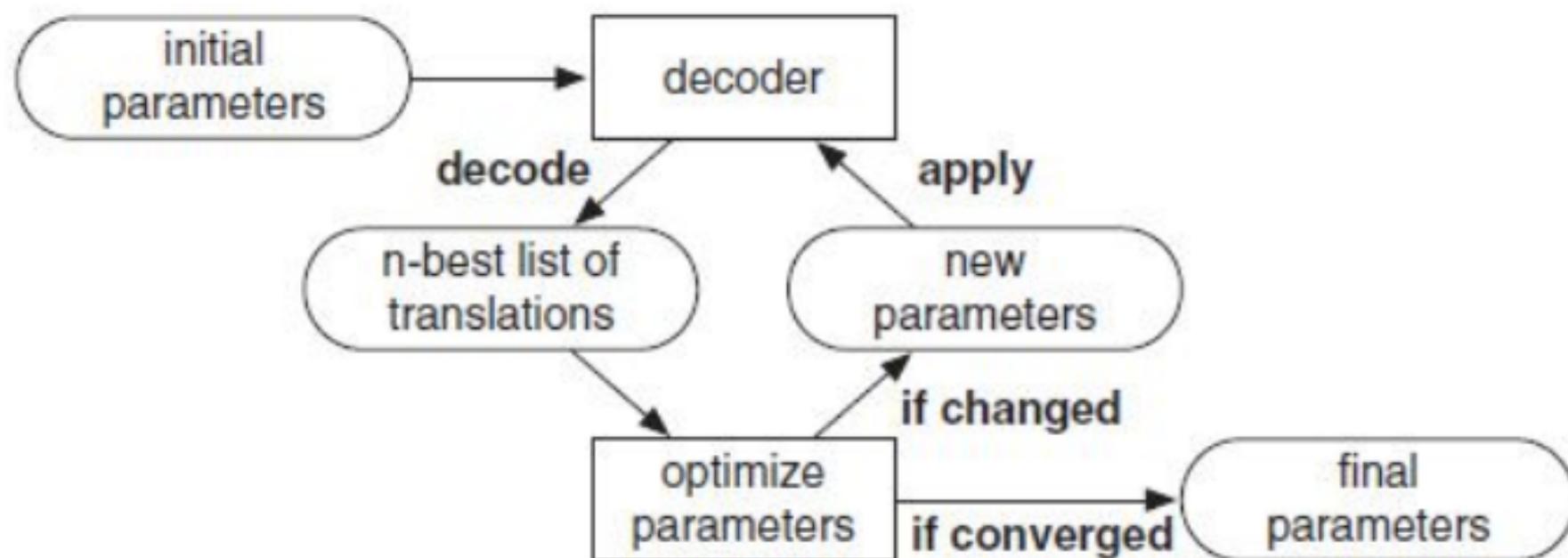
$$\ell(\hat{e}, \mathcal{E}) \mapsto [0, 1]$$

- Optimize the weight vector by making reference to this function

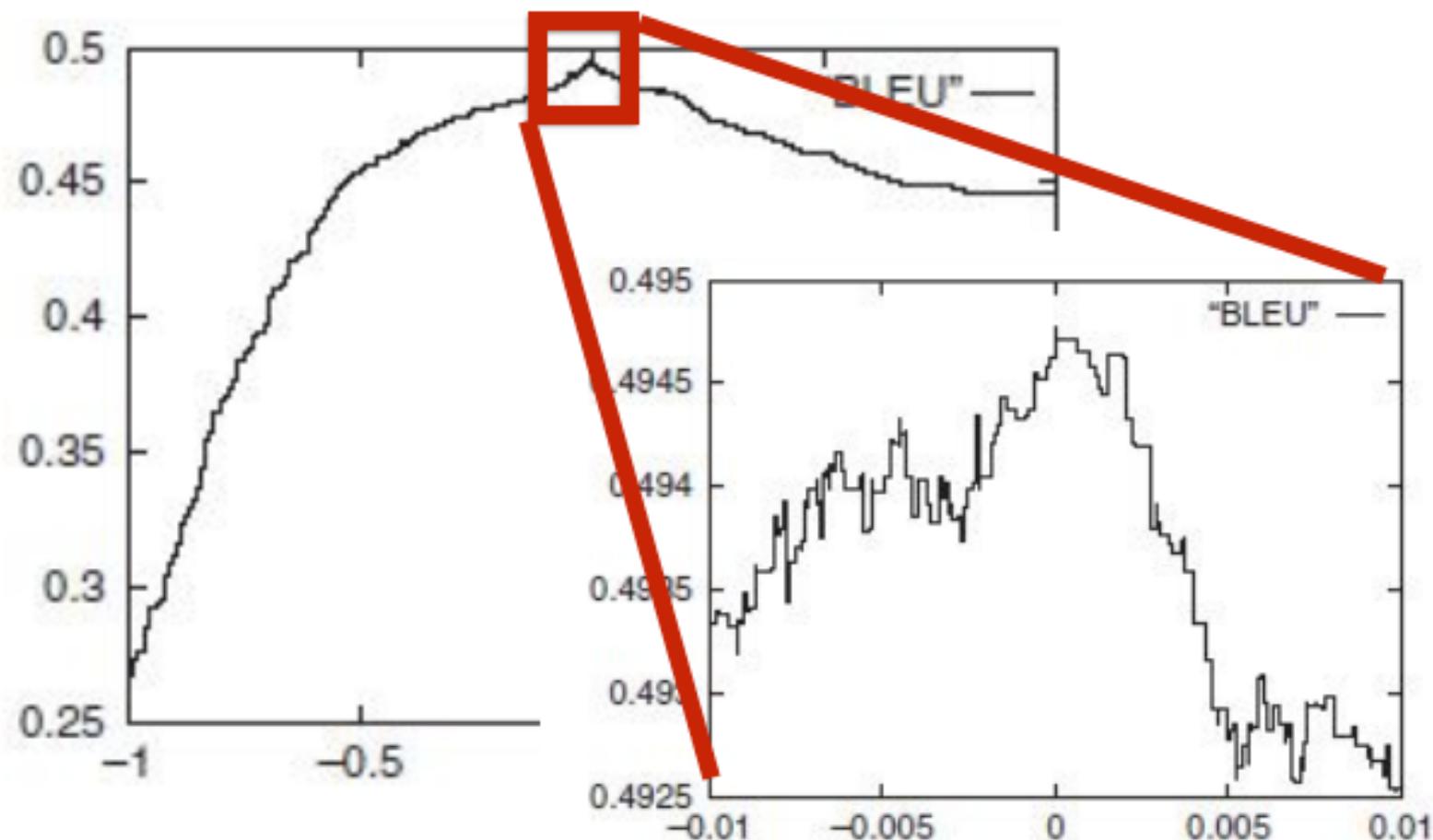
# MERT

- Minimum Error Rate Training
- Directly optimize for an automatic evaluation metric instead of likelihood
- Maximize the BLEU score on a **held out development set**
- Iteratively update the parameters by re-scoring n-best lists and comparing the highest scoring translation to the reference

# Iterative parameter tuning



# The effect on BLEU varying one parameter



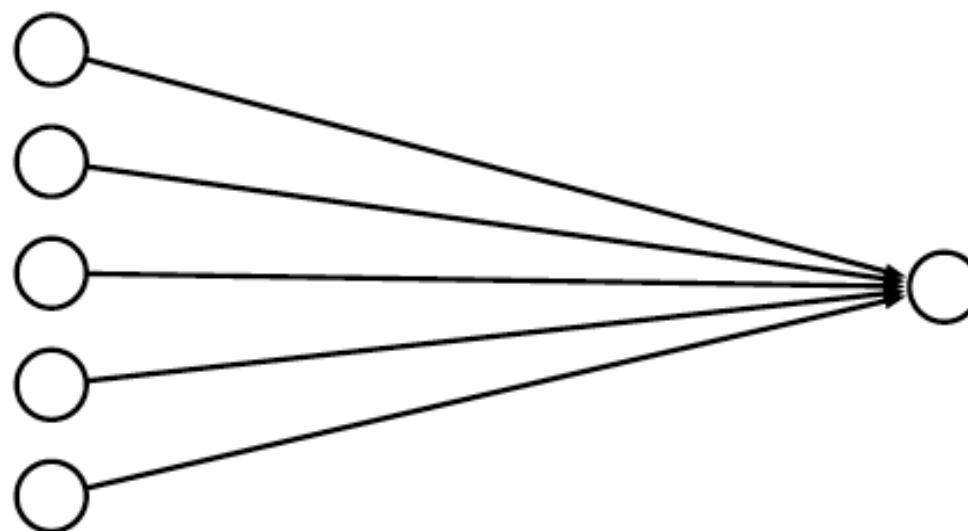
# Neural Machine Translation

## Linear Models

- We used before weighted linear combination of feature values  $h_j$  and weights  $\lambda_j$

$$\text{score}(\lambda, \mathbf{d}_i) = \sum_j \lambda_j h_j(\mathbf{d}_i)$$

- Such models can be illustrated as a "network"

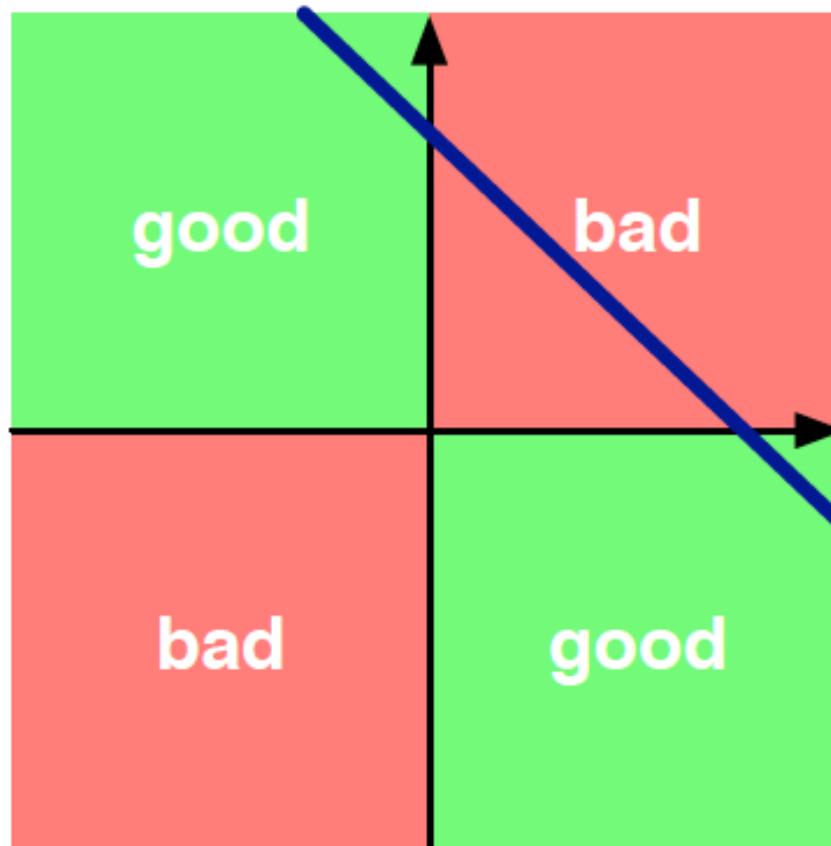


## Limits of Linearity

- We can give each feature a weight
- But not more complex value relationships, e.g,
  - any value in the range [0;5] is equally good
  - values over 8 are bad
  - higher than 10 is not worse

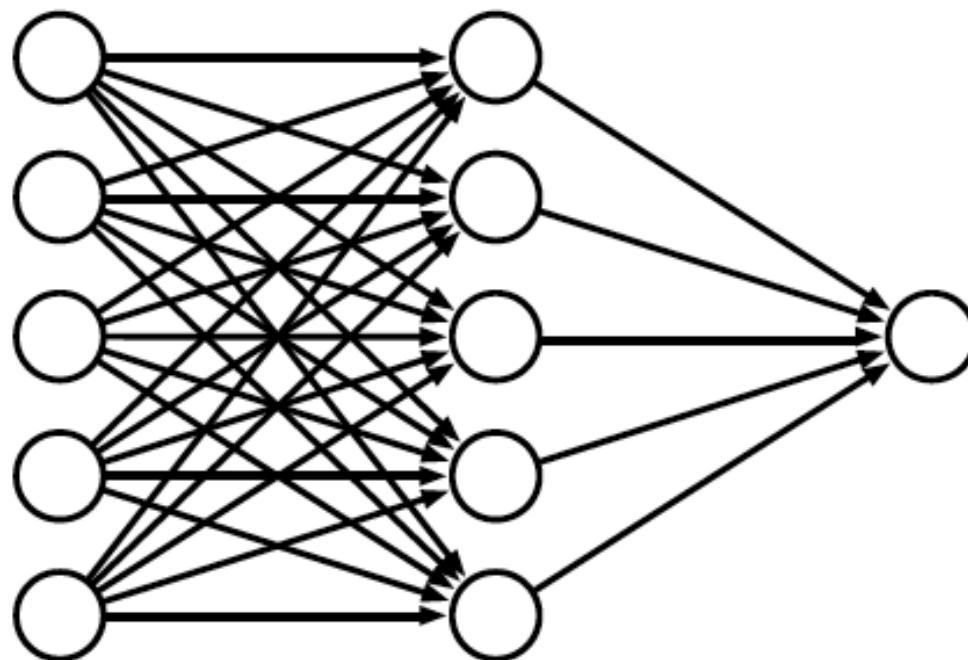
# XOR

- Linear models cannot model XOR



## Multiple Layers

- Add an intermediate ("hidden") layer of processing  
(each arrow is a weight)



- Have we gained anything so far?

# Non-Linearity

- Instead of computing a linear combination

$$\text{score}(\lambda, \mathbf{d}_i) = \sum_j \lambda_j h_j(\mathbf{d}_i)$$

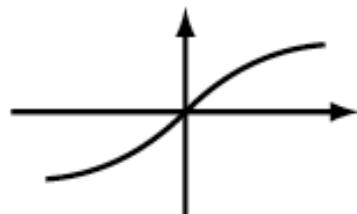
- Add a non-linear function

$$\text{score}(\lambda, \mathbf{d}_i) = f\left(\sum_j \lambda_j h_j(\mathbf{d}_i)\right)$$

- Popular choices

$$\tanh(x)$$

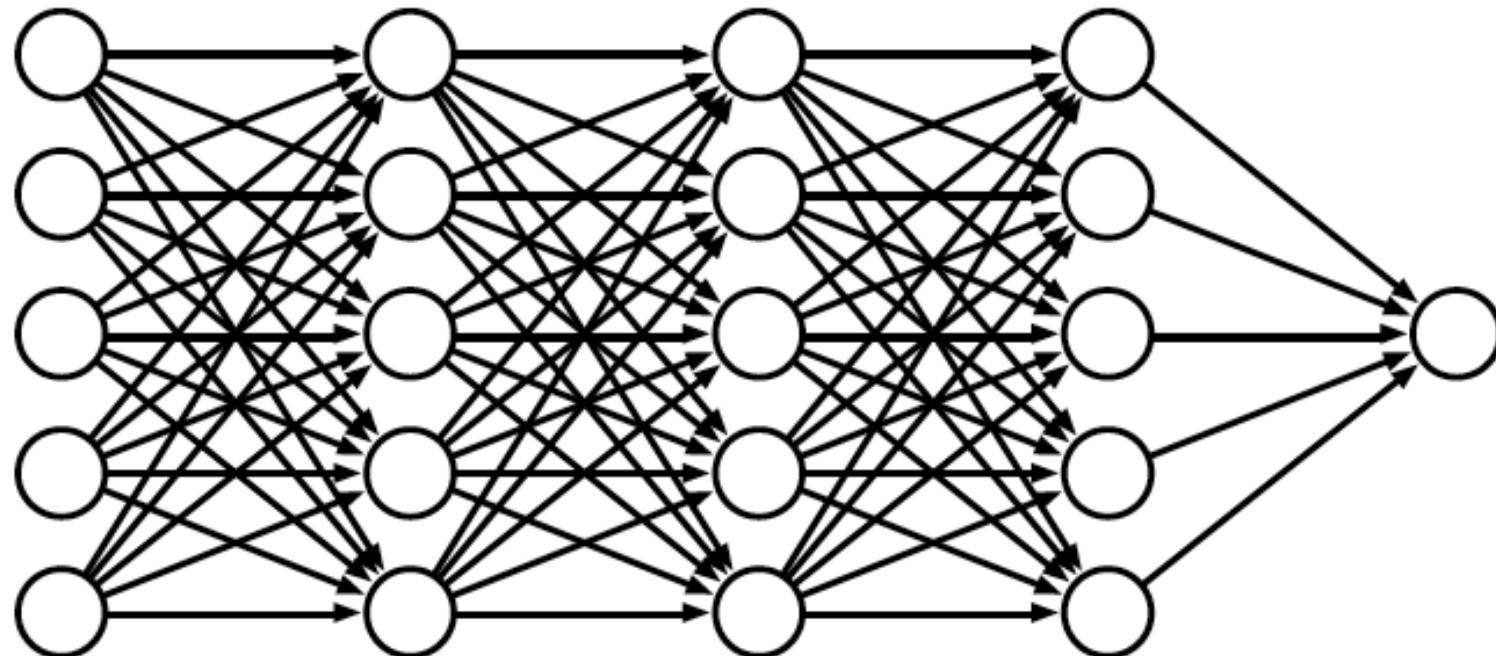
$$\text{sigmoid}(x) = \frac{1}{1+e^{-x}}$$



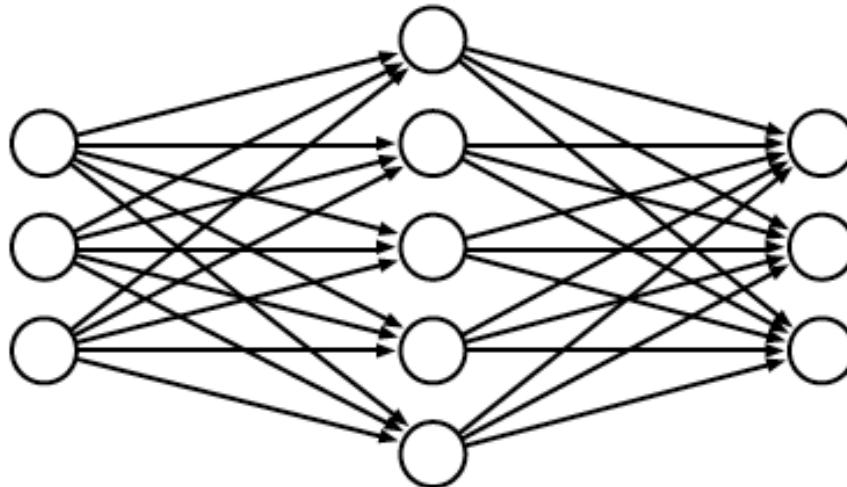
(sigmoid is also called the "logistic function")

# Deep Learning

- More layers = deep learning



## Neural Networks for Classification

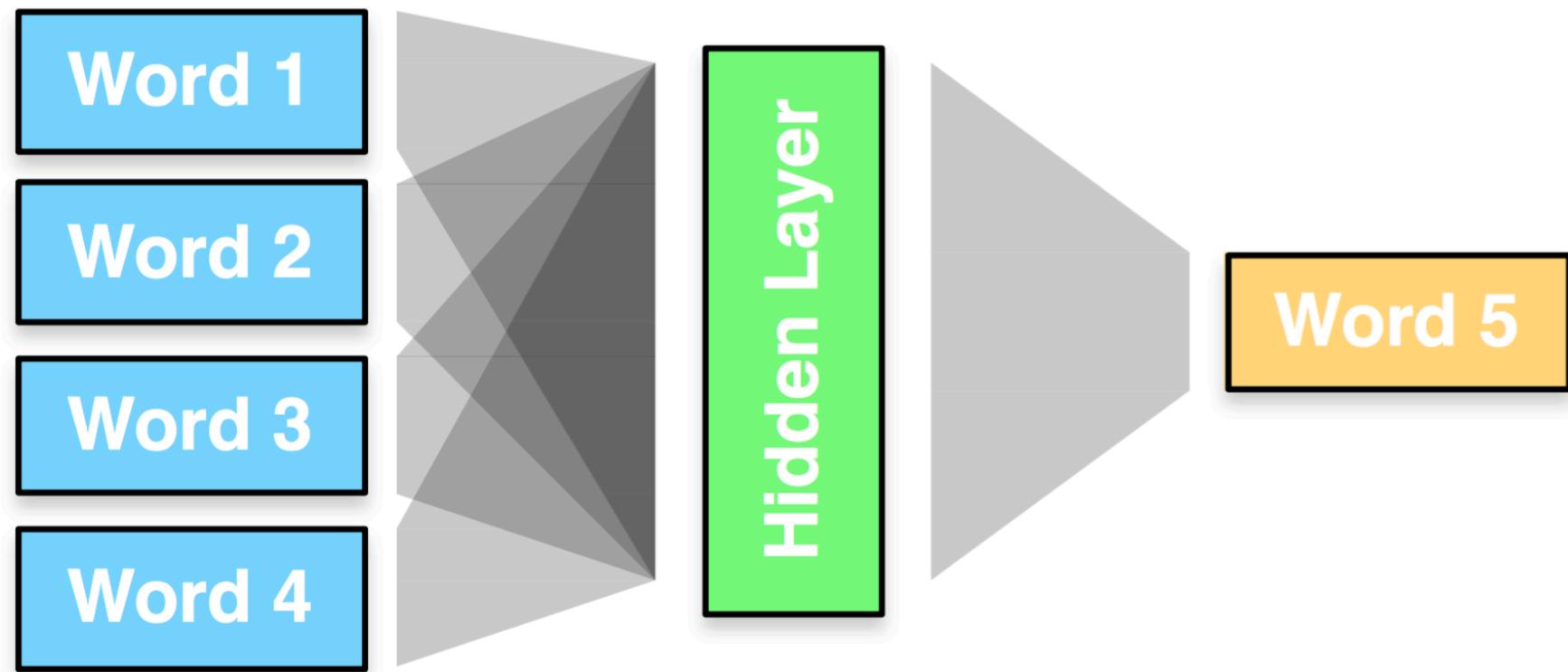


- Predict class: one output node per class
- Training data output: "One-hot vector", e.g.,  $\vec{y} = (0, 0, 1)^T$
- Prediction
  - predicted class is output node  $y_i$  with highest value
  - obtain posterior probability distribution by soft-max

$$\text{softmax}(y_i) = \frac{e^{y_i}}{\sum_j e^{y_j}}$$

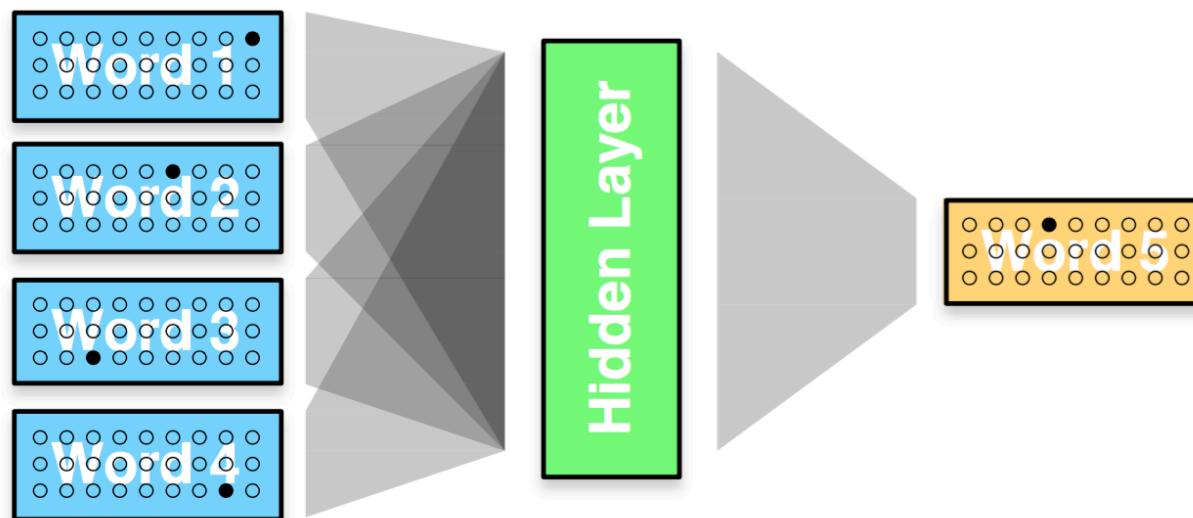
# Neural N-gram Language Models

$$p(\mathbf{e}) = \prod_{i=1}^{|\mathbf{e}|} p(e_i \mid e_{i-n+1}, \dots, e_{i-1})$$

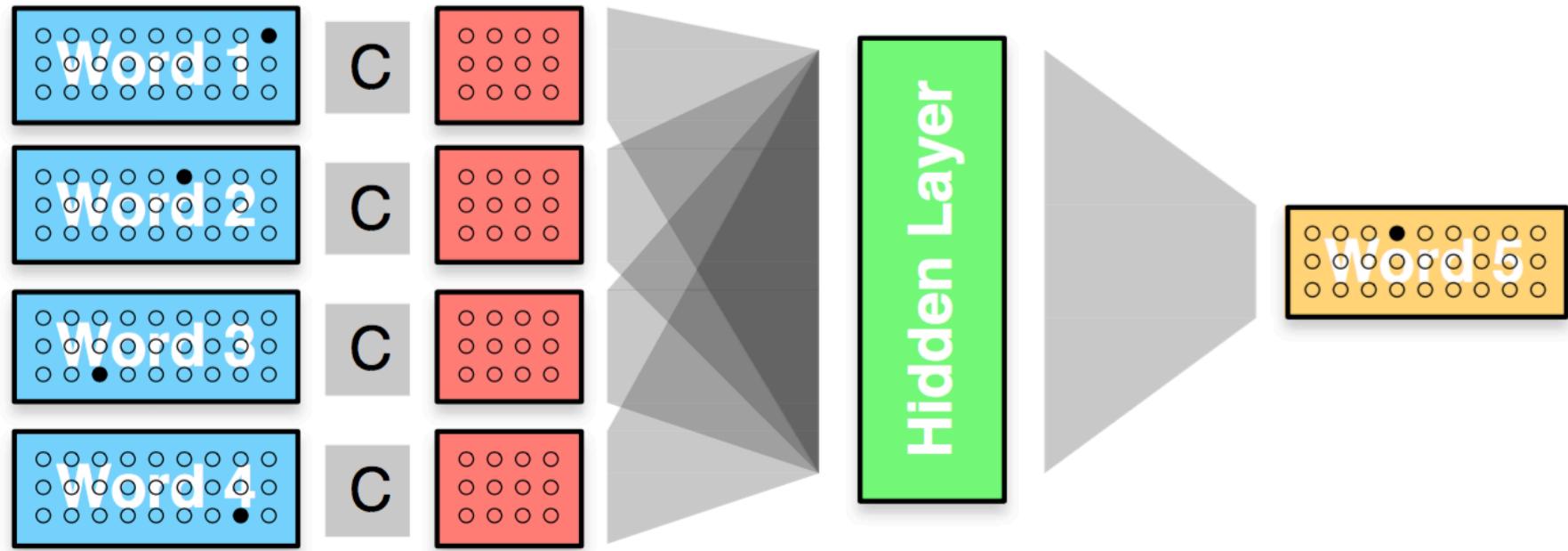


# Word Representation

- Words are represented with a one-hot vector, e.g.,
  - **dog** =  $(0,0,0,0,1,0,0,0,0,\dots)$
  - **cat** =  $(0,0,0,0,0,0,0,1,0,\dots)$
  - **eat** =  $(0,1,0,0,0,0,0,0,0,\dots)$
- That's a large vector!
- Remedies
  - limit to, say, 20,000 most frequent words, rest are OTHER

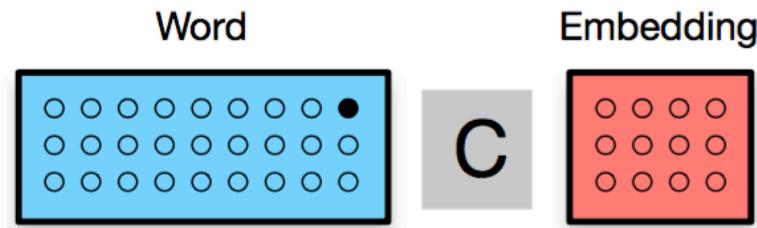


# Word embeddings



- Map each word first into a lower-dimensional real-valued space
- Shared weight matrix  $C$

# Word embeddings

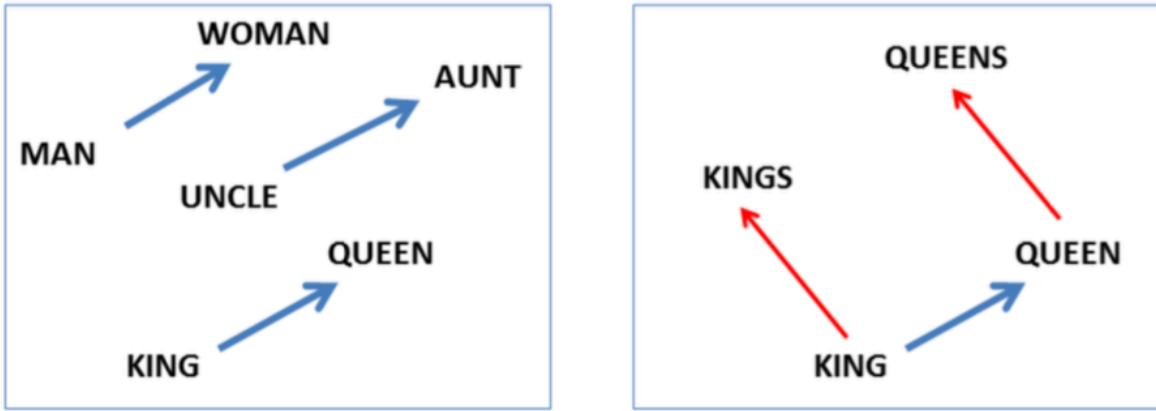


- By-product: embedding of word into continuous space
- Similar contexts → similar embedding
- Recall: distributional semantics

# Word embeddings



# Word embeddings

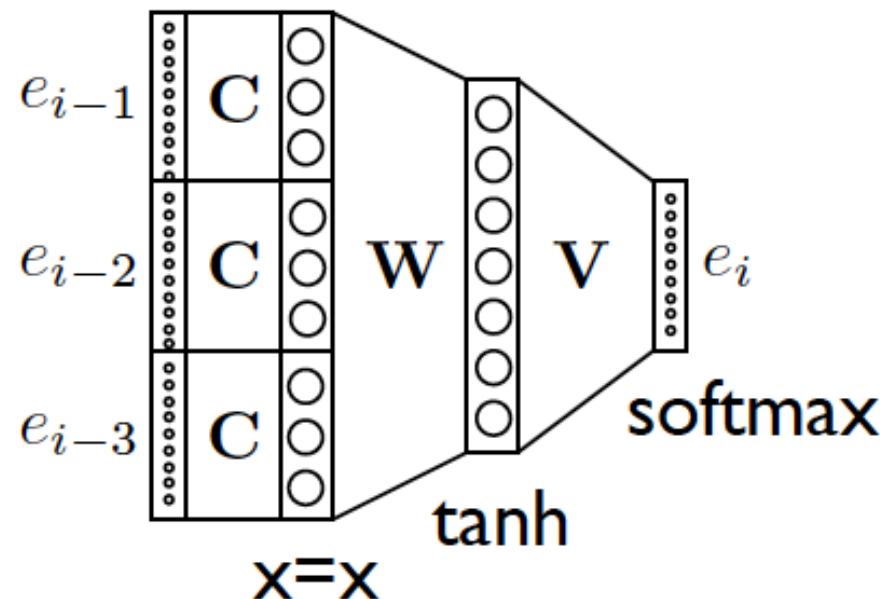


- Morphosyntactic regularities (Mikolov et al., 2013)
  - adjectives base form vs. comparative, e.g., **good**, **better**
  - nouns singular vs. plural, e.g., **year**, **years**
  - verbs present tense vs. past tense, e.g., **see**, **saw**
- Semantic regularities
  - **clothing** is to **shirt** as **dish** is to **bowl**
  - evaluated on human judgment data of semantic similarities

# Bengio et al. (2003)

$$p(\mathbf{e}) = \prod_{i=1}^{|\mathbf{e}|} p(e_i \mid e_{i-n+1}, \dots, e_{i-1})$$

$$p(e_i \mid e_{i-n+1}, \dots, e_{i-1}) =$$

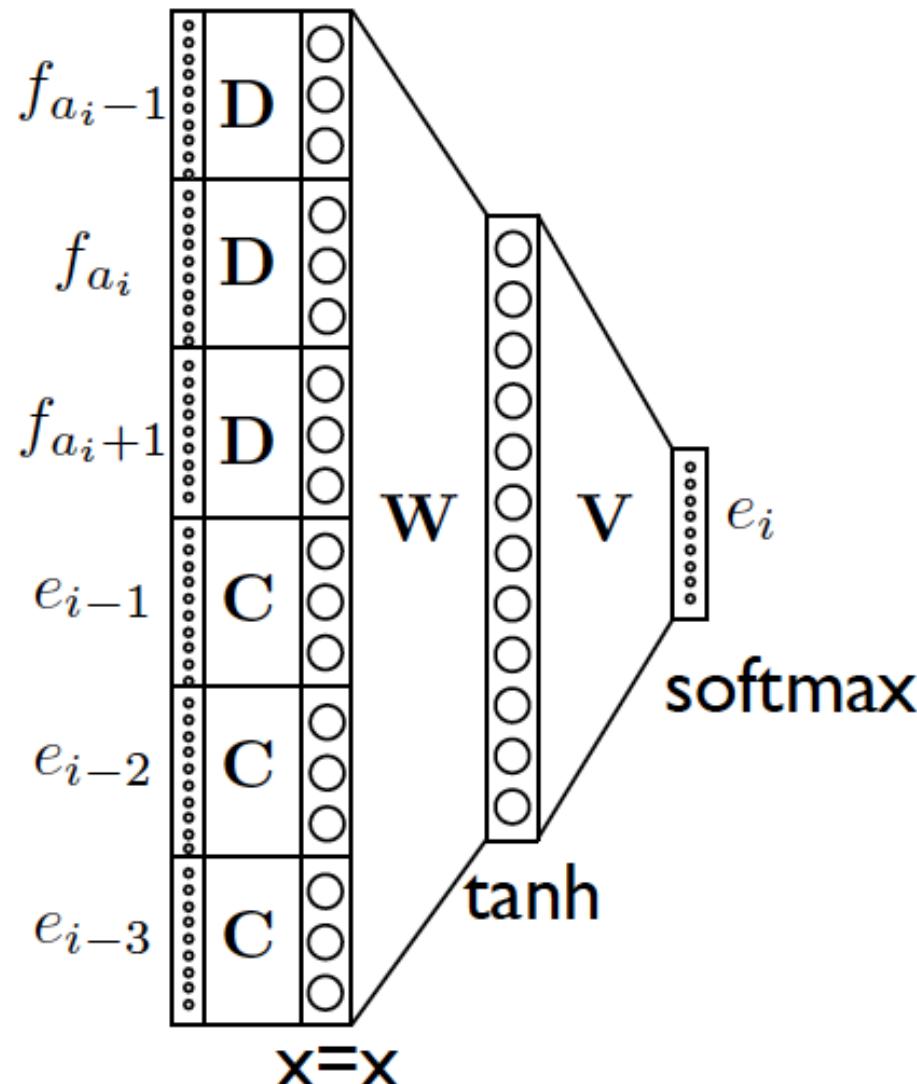


# Devlin et al. (2014)

- Turn Bengio et al. (2003) into a translation model
- Conditional model; generate the next English word conditioned on
  - The previous  $n$  English words you generated
  - The aligned source word, and its  $m$  neighbors

$$p(\mathbf{e} \mid \mathbf{f}, \mathbf{a}) = \prod_{i=1}^{|e|} p(e_i \mid e_{i-2}, e_{i-1}, f_{a_i-1}, f_{a_i}, f_{a_i+1})$$

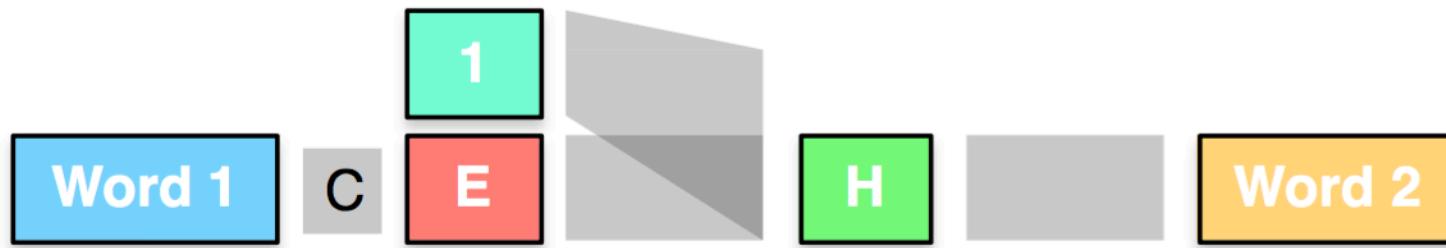
$$p(e_i \mid e_{i-2}, e_{i-1}, f_{a_i-1}, f_{a_i}, f_{a_i+1}) =$$



# Summary

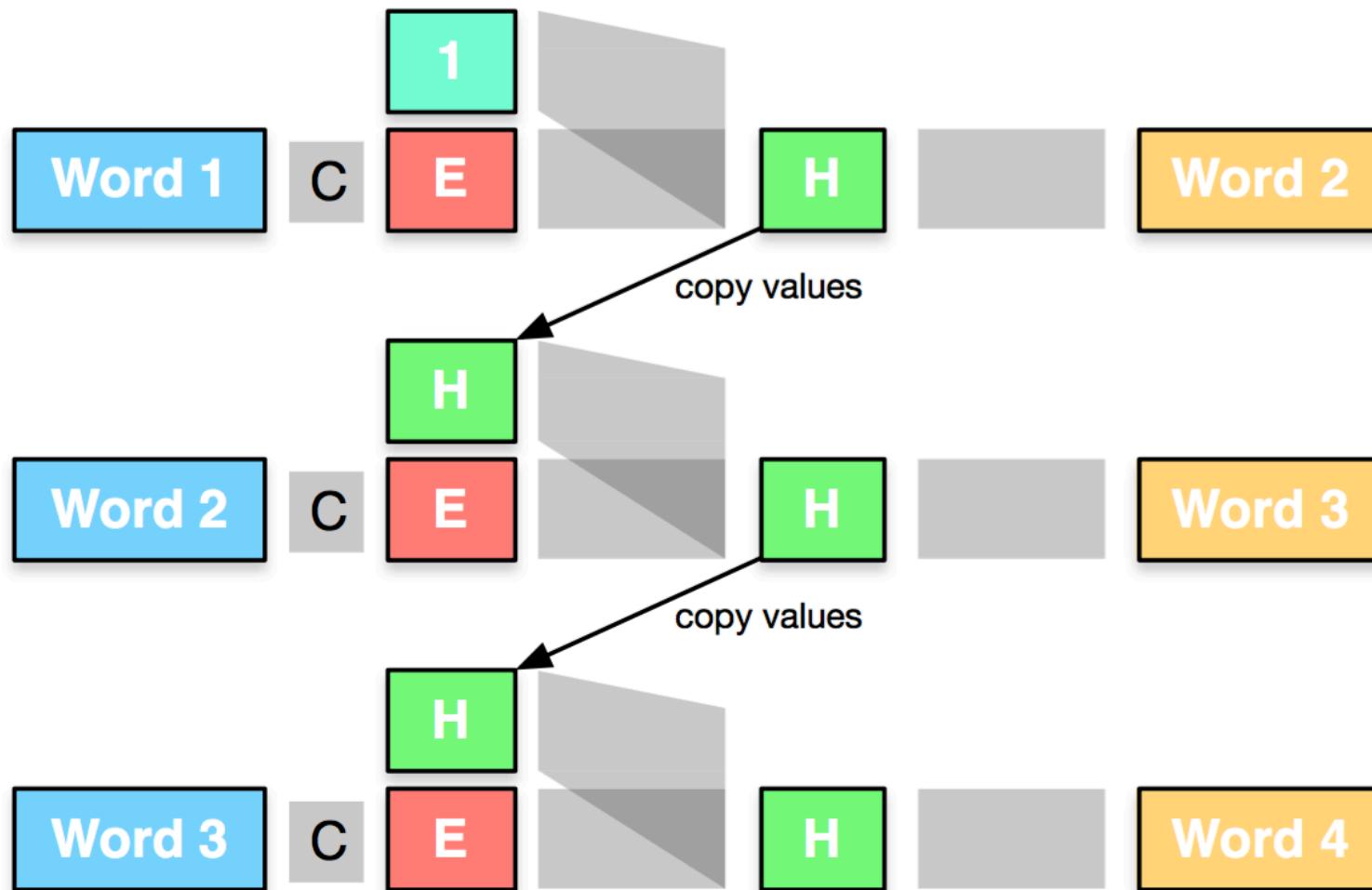
- Two problems in standard statistical models
  - We don't condition on enough stuff
  - We don't know what features to use when we condition on lots of structure
- **Punchline: Neural networks let us condition on a lot of stuff without an exponential growth in parameters**

# Recurrent NN

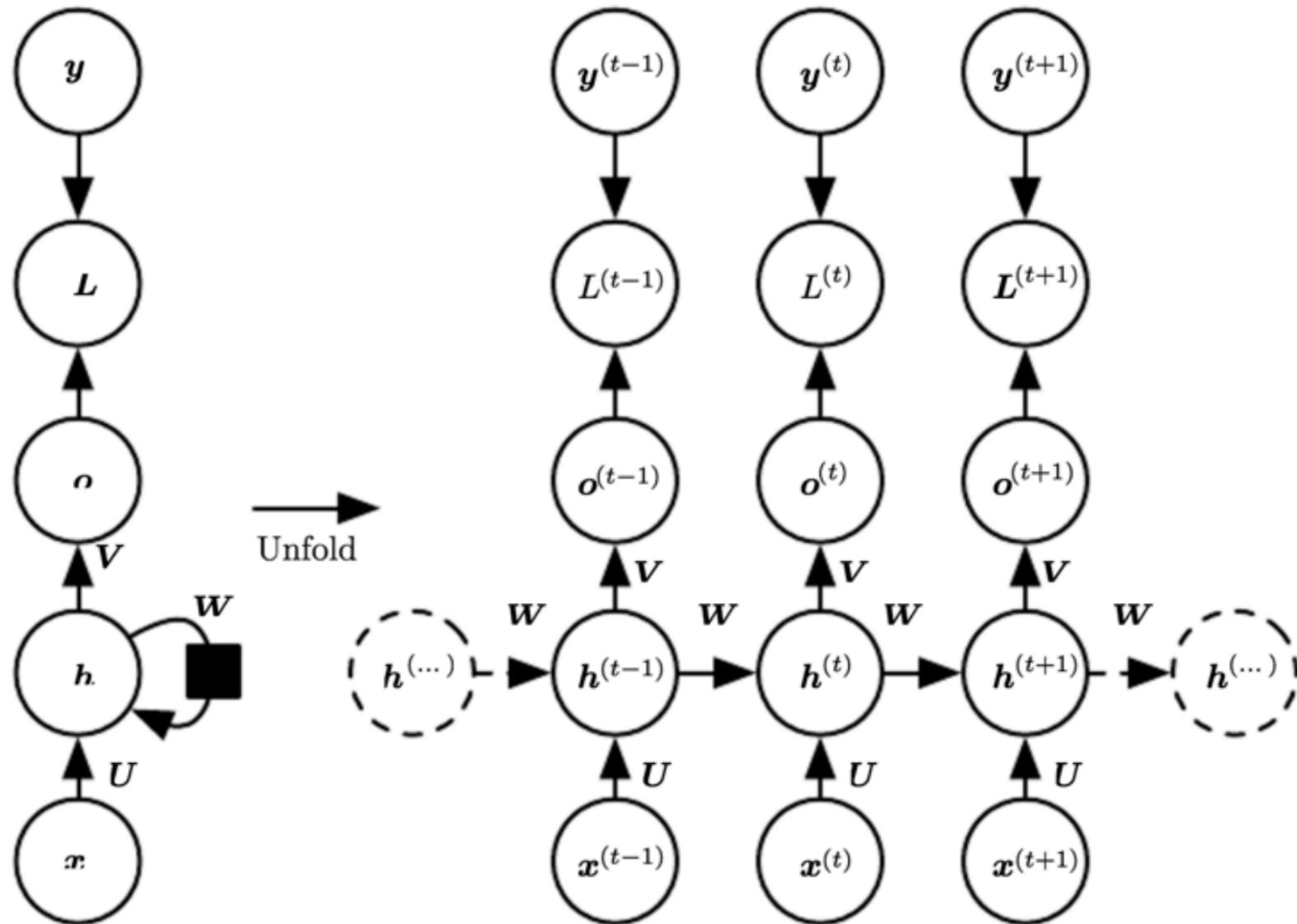


- Start: predict second word from first
- Mystery layer with nodes all with value 1

# Recurrent NN



# Recurrent NN

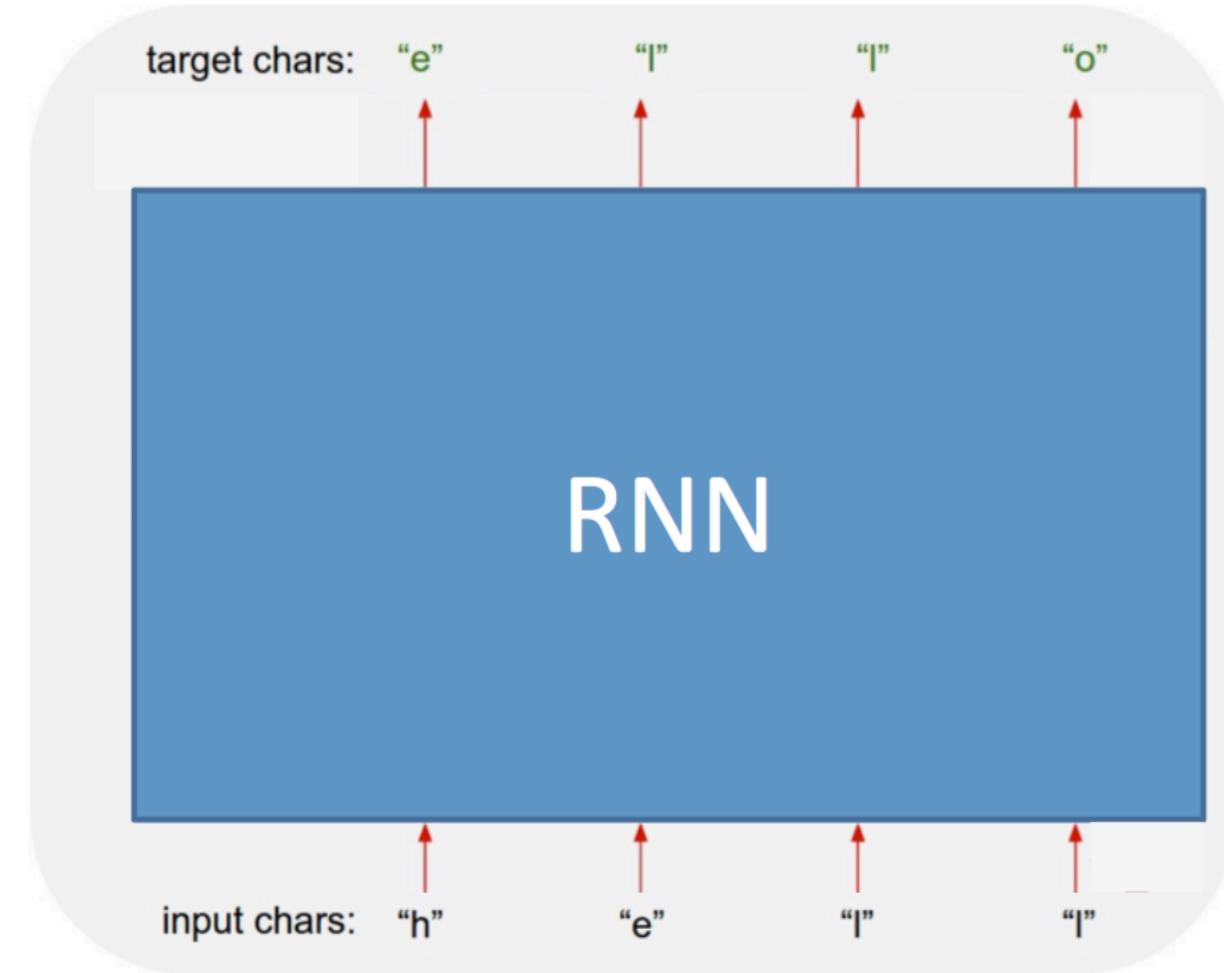


# RNN Language Model

**Character-level  
language model  
example**

Vocabulary:  
[h,e,l,o]

Example training  
sequence:  
**“hello”**

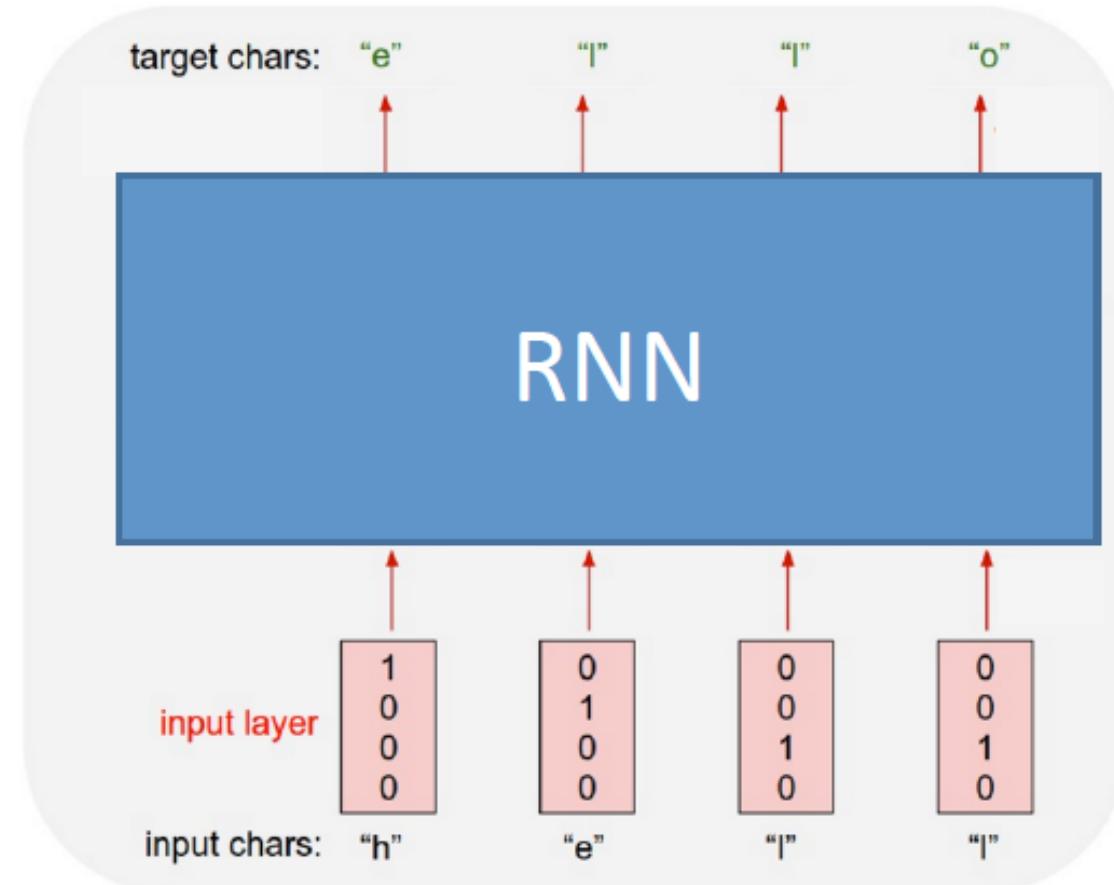


# RNN – *Input Layer*

**Character-level  
language model  
example**

Vocabulary:  
[h,e,l,o]

Example training  
sequence:  
**“hello”**



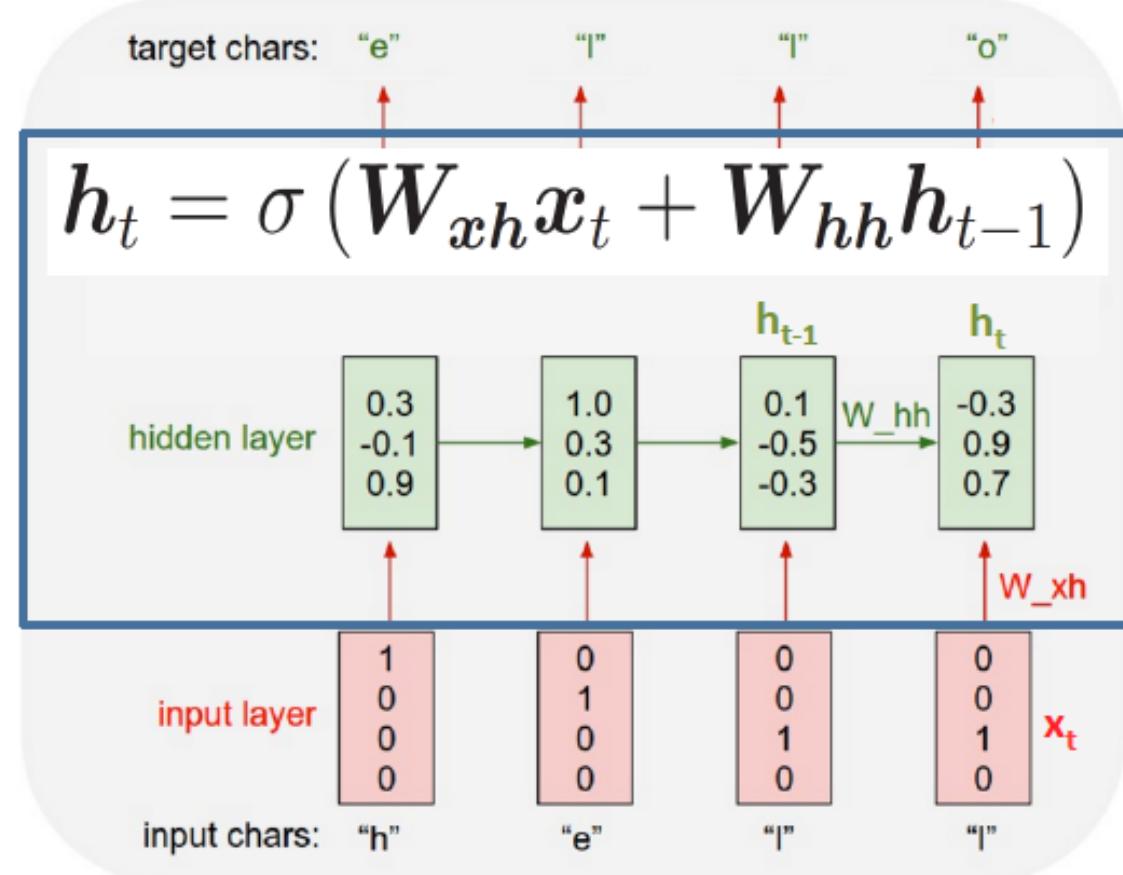
(Picture adapted from Andrej Karpathy)

# RNN – Hidden Layer

**Character-level  
language model  
example**

Vocabulary:  
[h,e,l,o]

Example training  
sequence:  
**“hello”**



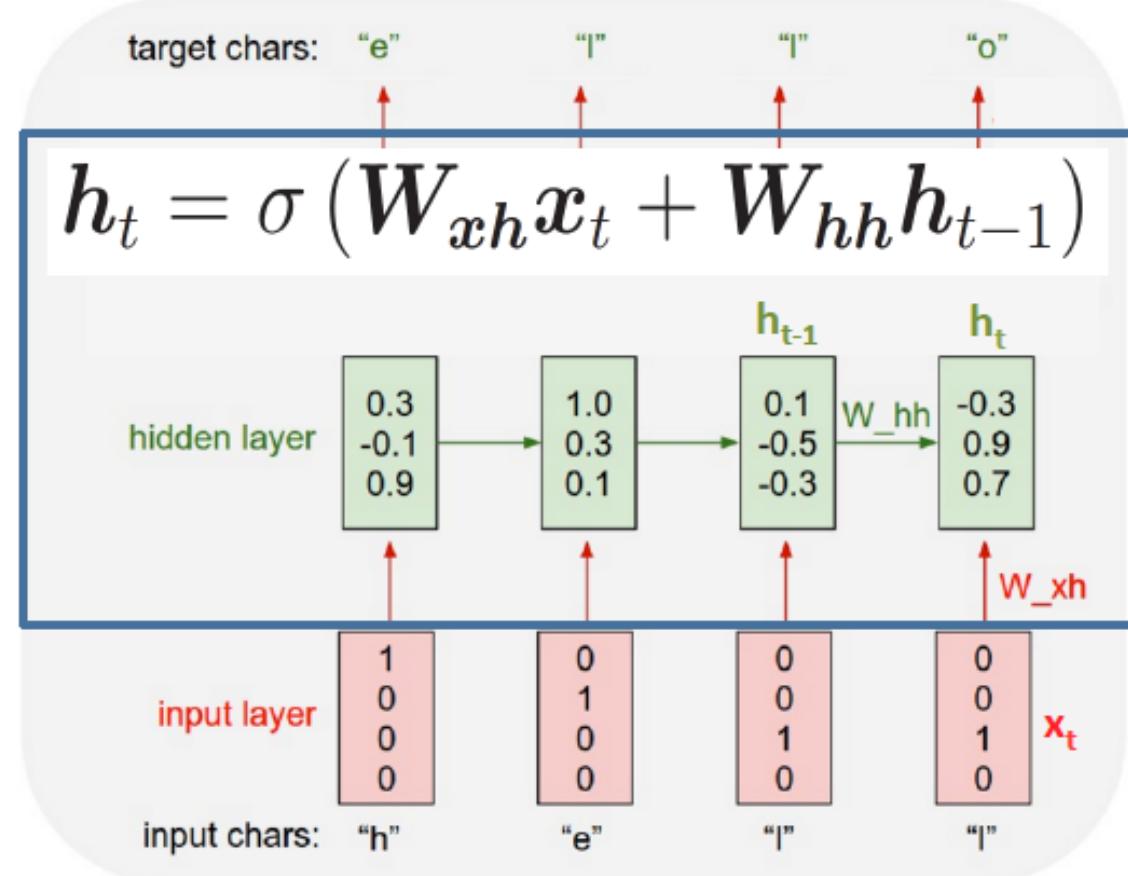
(Picture adapted from Andrej Karpathy)

# RNN – *Hidden Layer*

**Character-level  
language model  
example**

Vocabulary:  
[h,e,l,o]

Example training  
sequence:  
**“hello”**



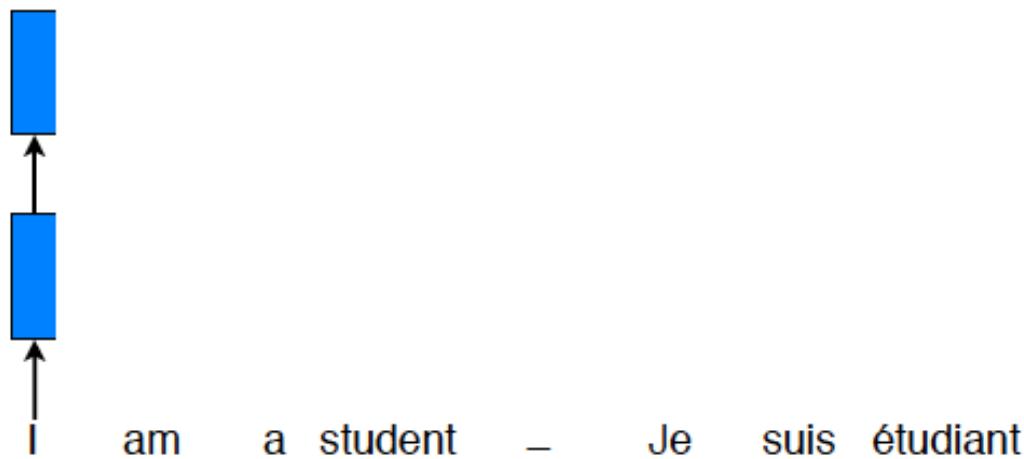
**RNNs to represent sequences!**

# Neural Machine Translation (NMT)

I am a student – Je suis étudiant

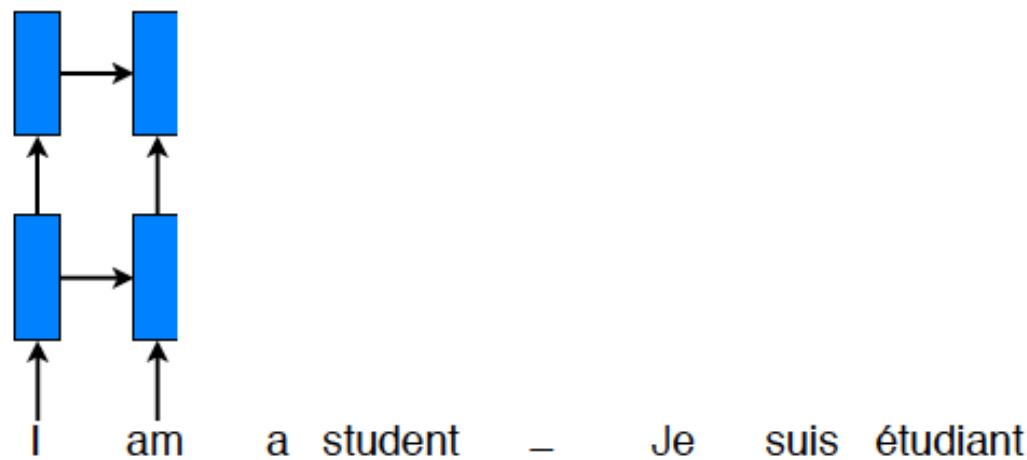
- Model  $P(\text{target} \mid \text{source})$  directly.

# Neural Machine Translation (NMT)



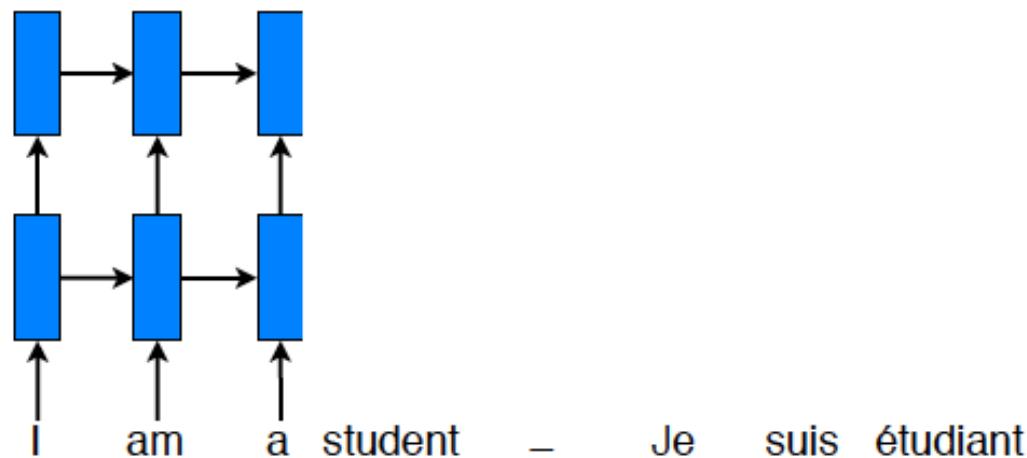
- RNNs trained **end-to-end** (Sutskever et al., 2014).

# Neural Machine Translation (NMT)



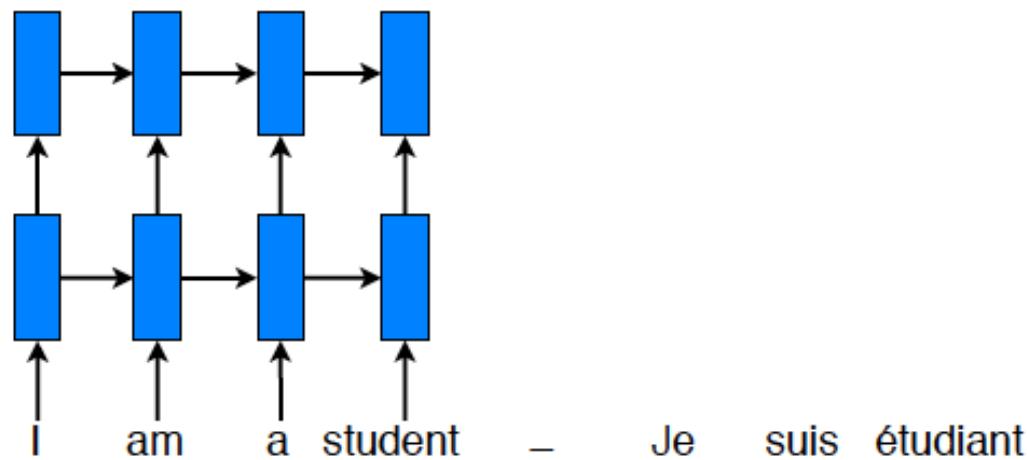
- RNNs trained **end-to-end** (Sutskever et al., 2014).

# Neural Machine Translation (NMT)



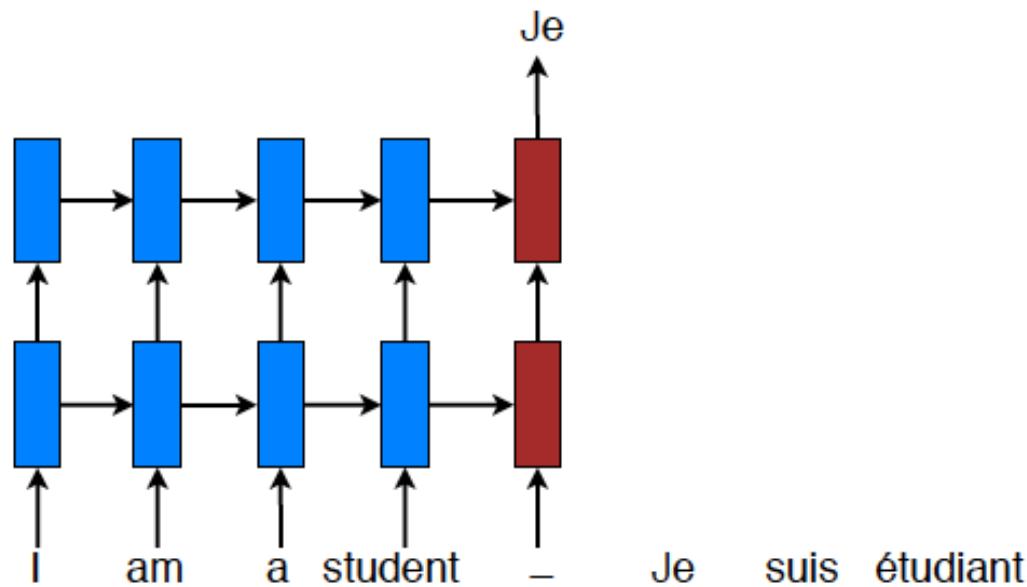
- RNNs trained **end-to-end** (Sutskever et al., 2014).

# Neural Machine Translation (NMT)



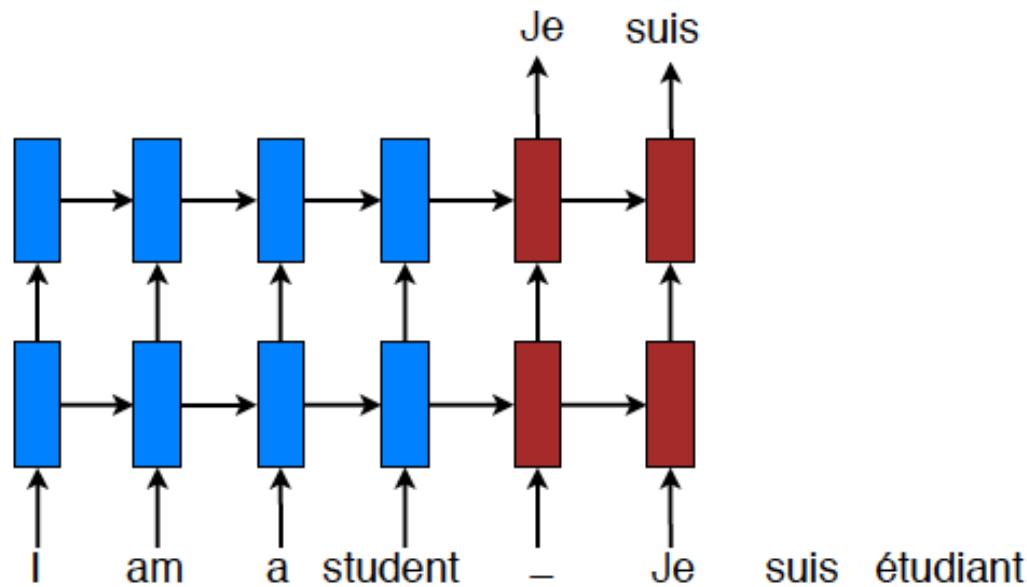
- RNNs trained **end-to-end** (Sutskever et al., 2014).

# Neural Machine Translation (NMT)



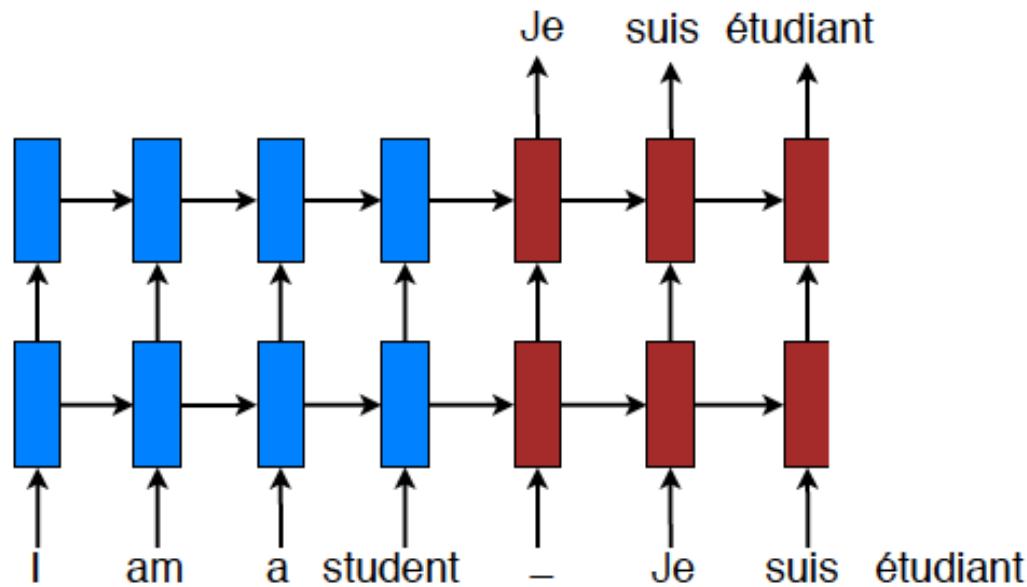
- RNNs trained **end-to-end** (Sutskever et al., 2014).

# Neural Machine Translation (NMT)



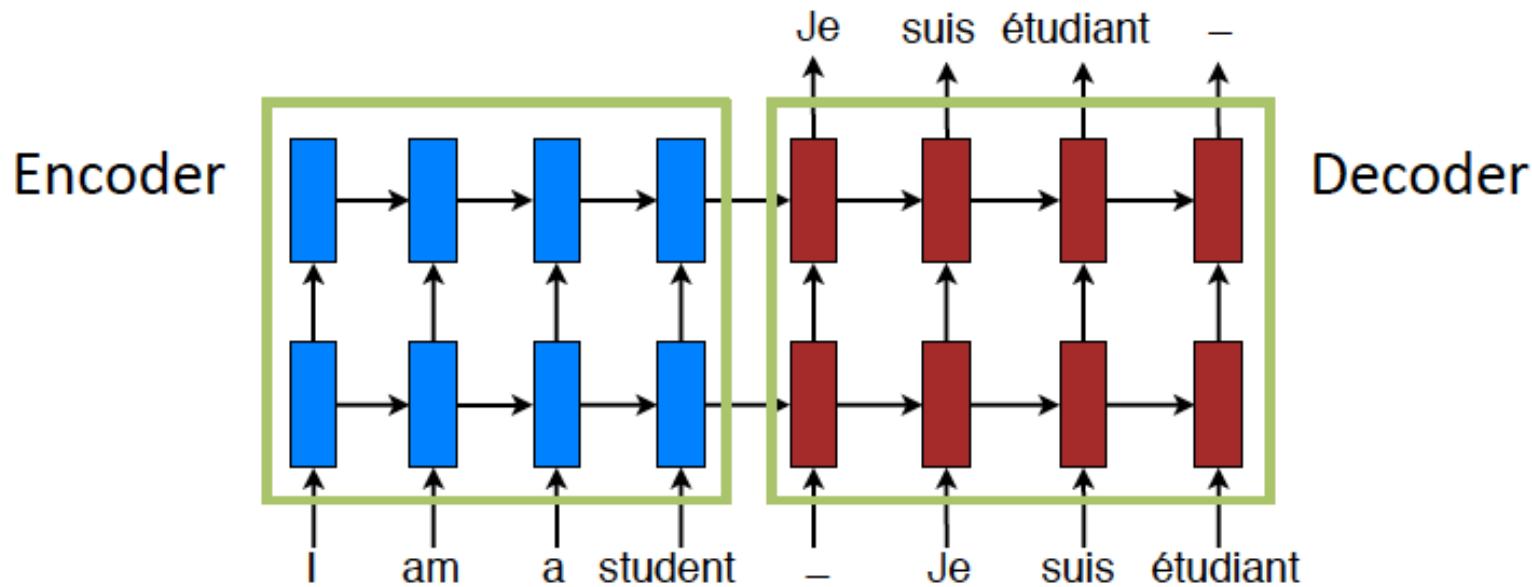
- RNNs trained end-to-end (Sutskever et al., 2014).

# Neural Machine Translation (NMT)



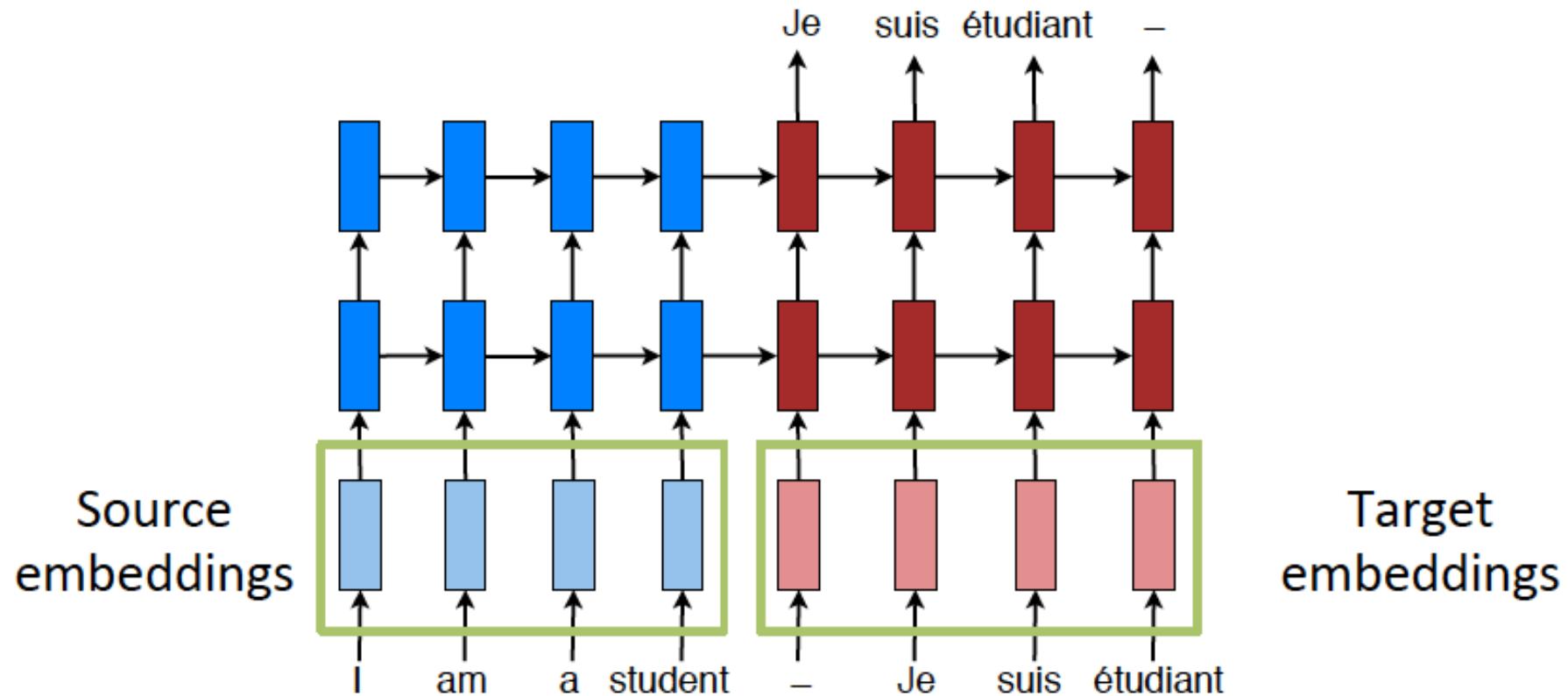
- RNNs trained end-to-end (Sutskever et al., 2014).

# Neural Machine Translation (NMT)



- RNNs trained **end-to-end** (Sutskever et al., 2014).
- Encoder-decoder approach.

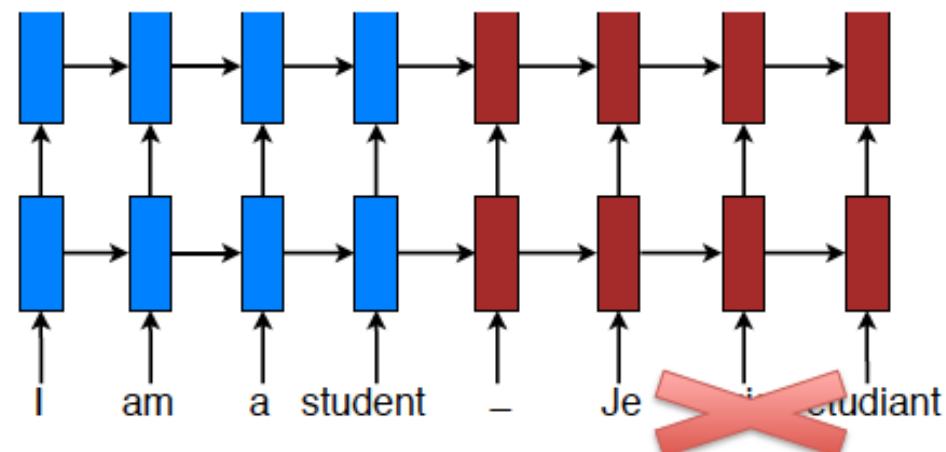
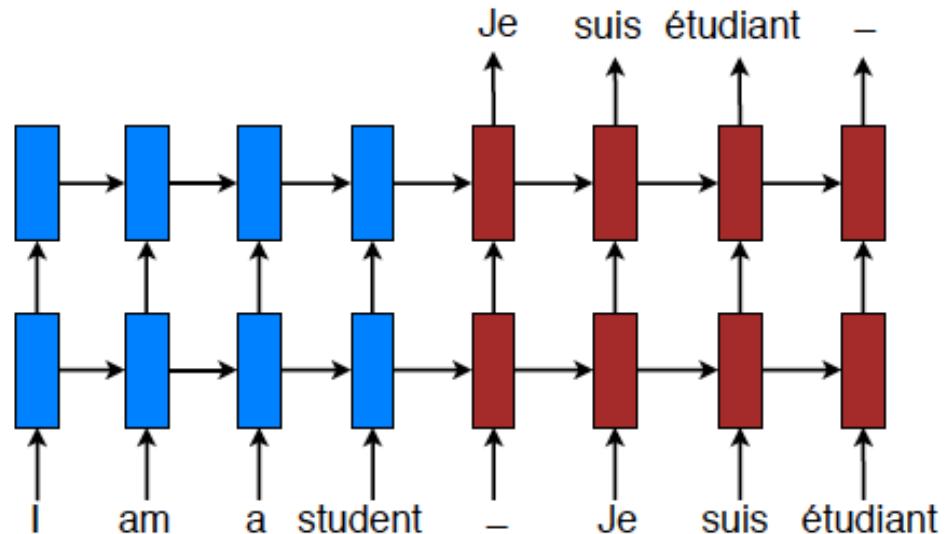
# Word Embeddings



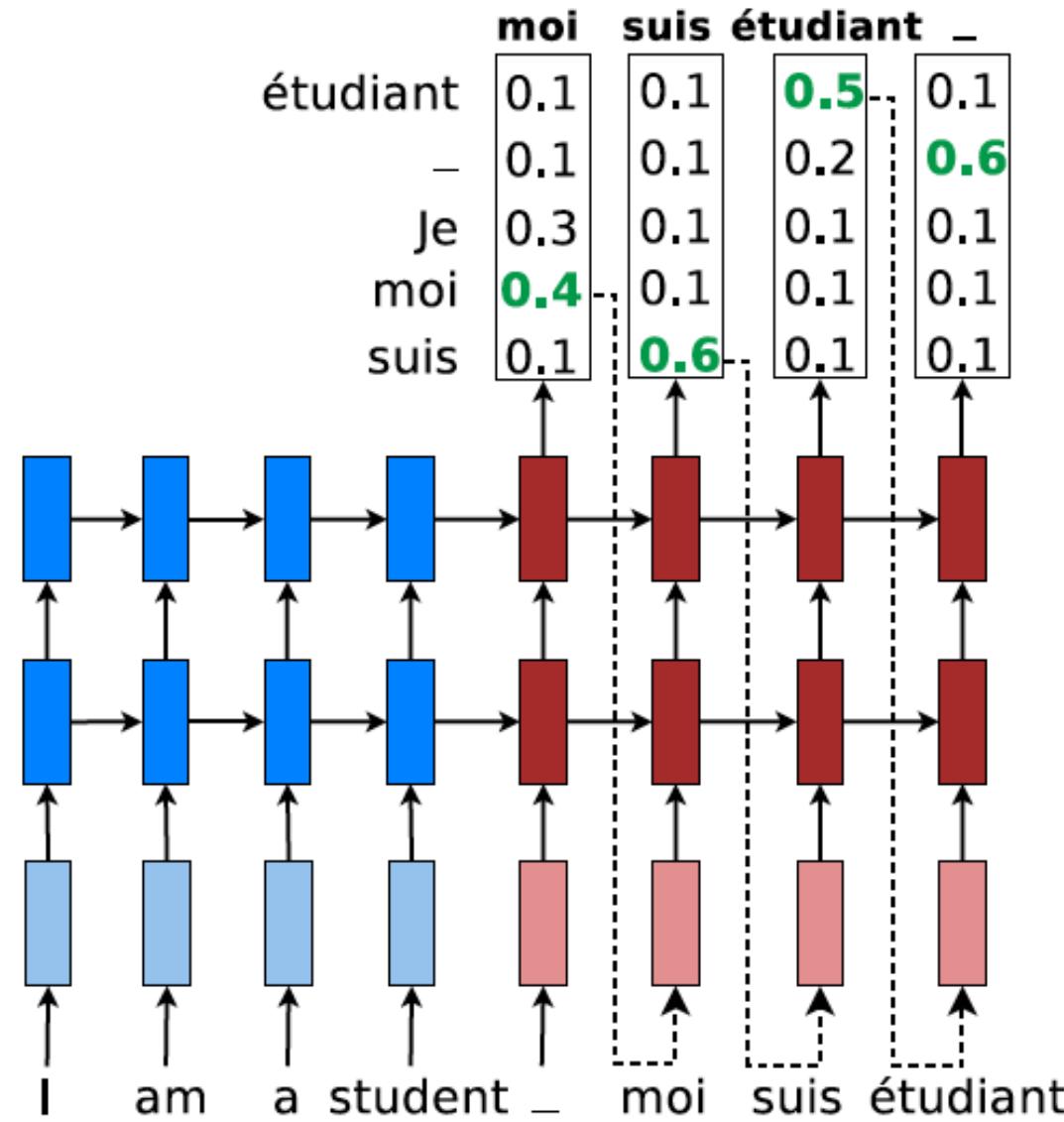
- Randomly initialized, one for each language.
  - Learnable parameters.

# Training vs. Testing

- *Training*
  - Correct translations are available.
- *Testing*
  - Only source sentences are given.



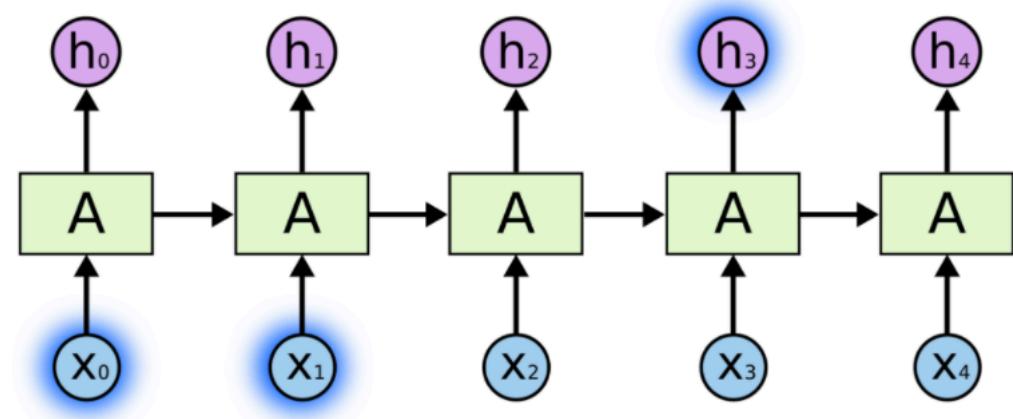
# Testing



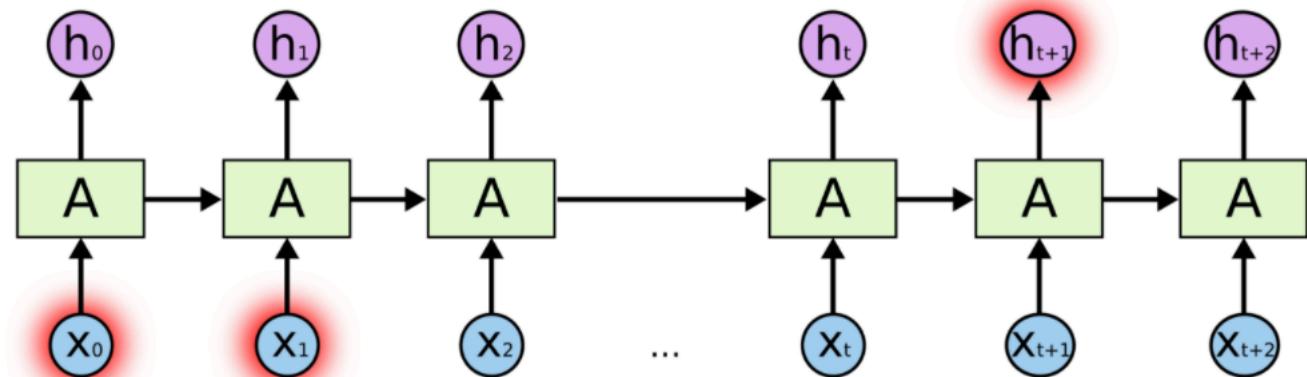
- Feed the most likely word

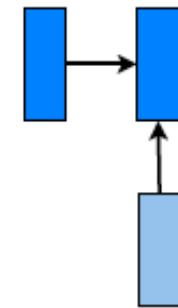
# The problem of long term dependencies

the clouds are in the sky

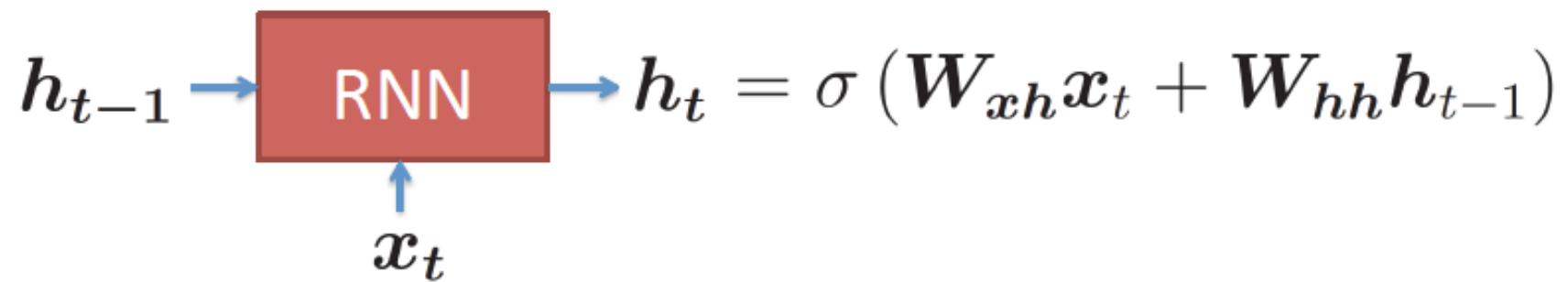


I grew up in France... I speak fluent French





## Recurrent types – vanilla RNN



Vanishing gradient problem!

# Vanishing gradients

$$\mathbf{h}_t = \sigma(\mathbf{W}_{xh}\mathbf{x}_t + \mathbf{W}_{hh}\mathbf{h}_{t-1})$$

$$\frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} = \text{diag}(\sigma'(\dots)) \mathbf{W}_{hh}^\top$$

Chain Rule

$$\left\| \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \right\| \leq \gamma \lambda_1$$

Bound Rules

$$\left\| \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-k}} \right\| \leq (\gamma \lambda_1)^k$$

Chain Rule

(Pascanu et al., 2013)

# Vanishing gradients

$$\mathbf{h}_t = \sigma(\mathbf{W}_{xh}\mathbf{x}_t + \mathbf{W}_{hh}\mathbf{h}_{t-1})$$

$$\frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} = \text{diag}(\sigma'(\dots)) \mathbf{W}_{hh}^\top$$

Chain Rule

$$\left\| \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \right\| \leq \gamma \lambda_1$$

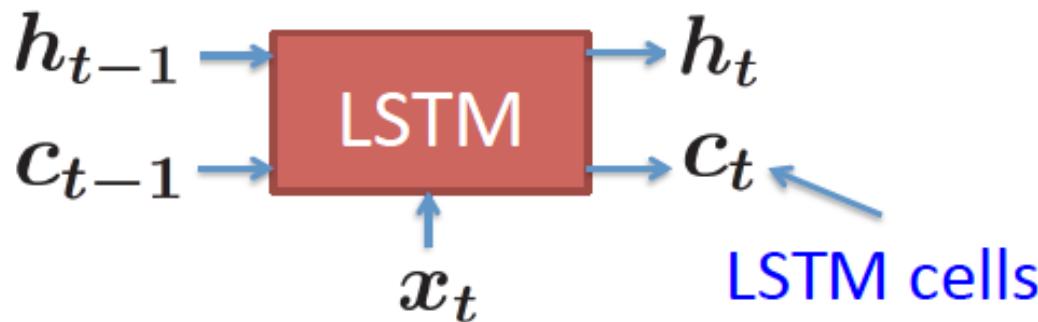
Bound Rules

$$\left\| \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-k}} \right\| \leq (\gamma \lambda_1)^k \rightarrow 0 \text{ if } \lambda_1 < \frac{1}{\gamma}$$

Sufficient Cond

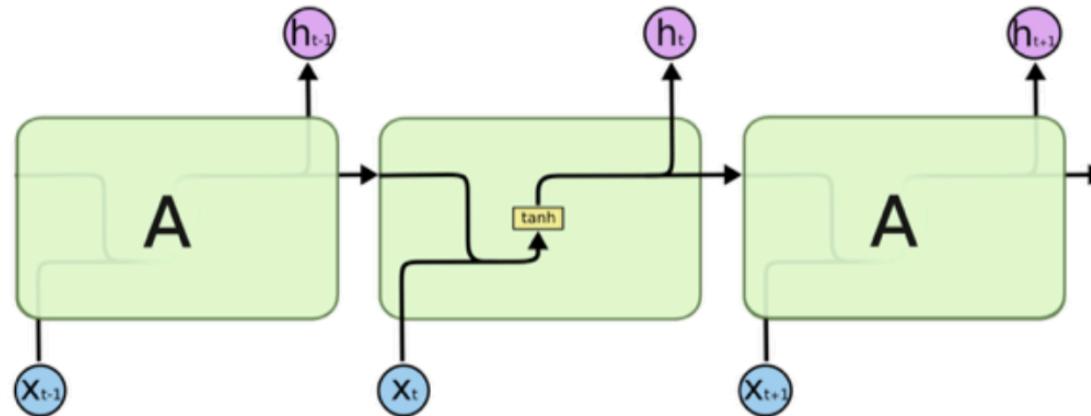
(Pascanu et al., 2013)

# Recurrent types – LSTM

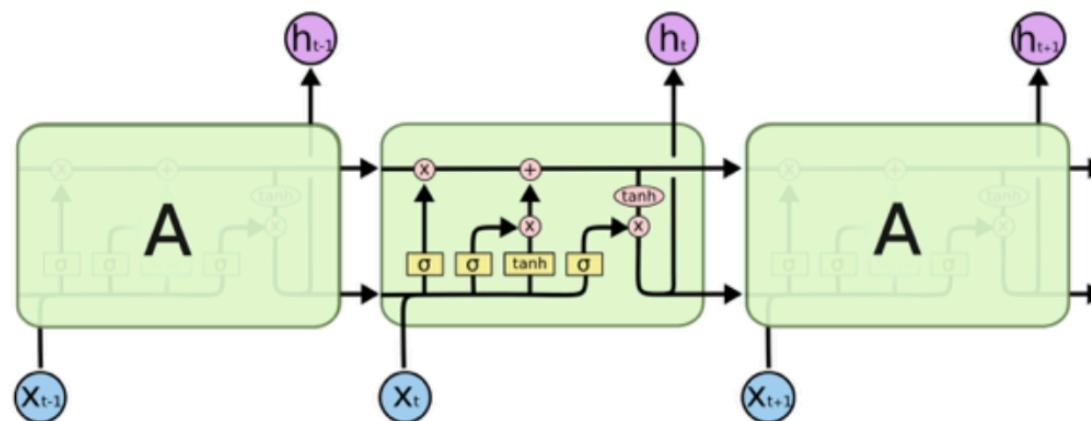


- Long-Short Term Memory (LSTM)
  - (Hochreiter & Schmidhuber, 1997)
- LSTM cells are **additively** updated
  - Make backprop through time easier.

RNN:



LSTM:



Neural Network Layer

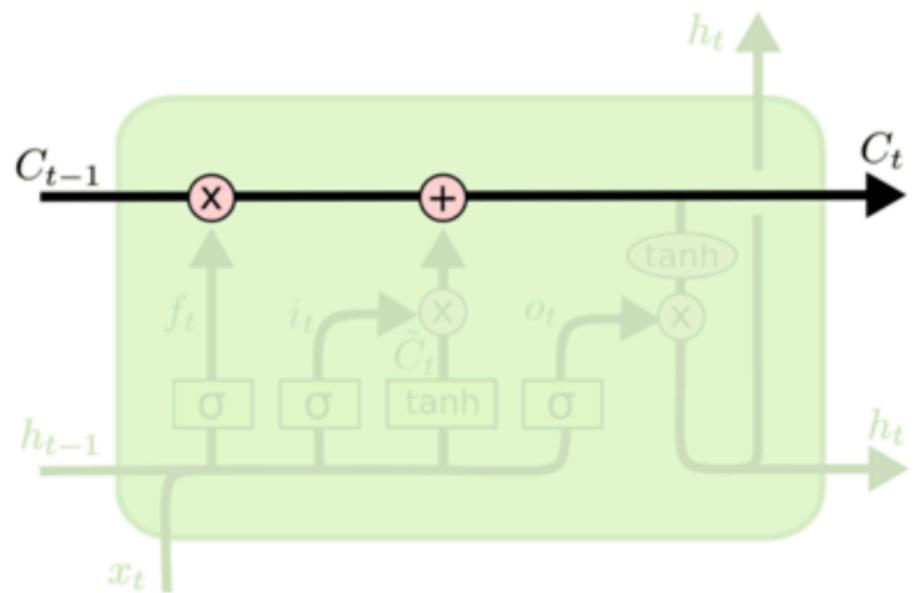
Pointwise Operation

Vector Transfer

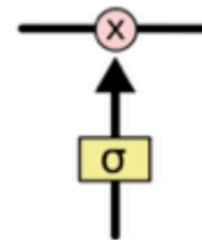
Concatenate

Copy

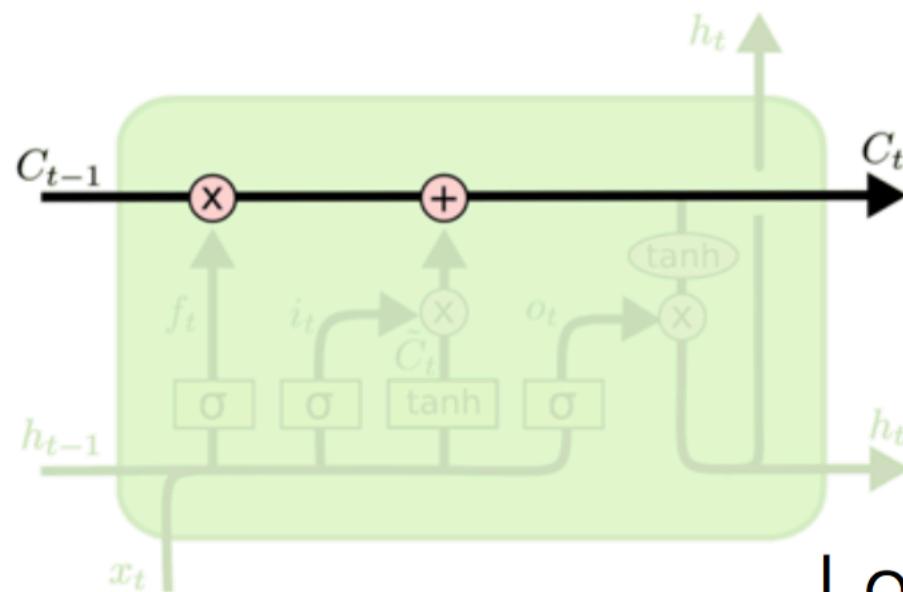
“The cell”



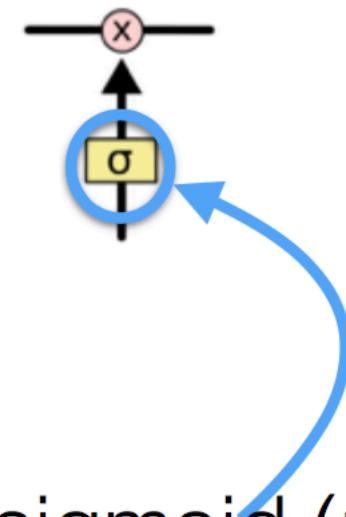
A Gate:



“The cell”

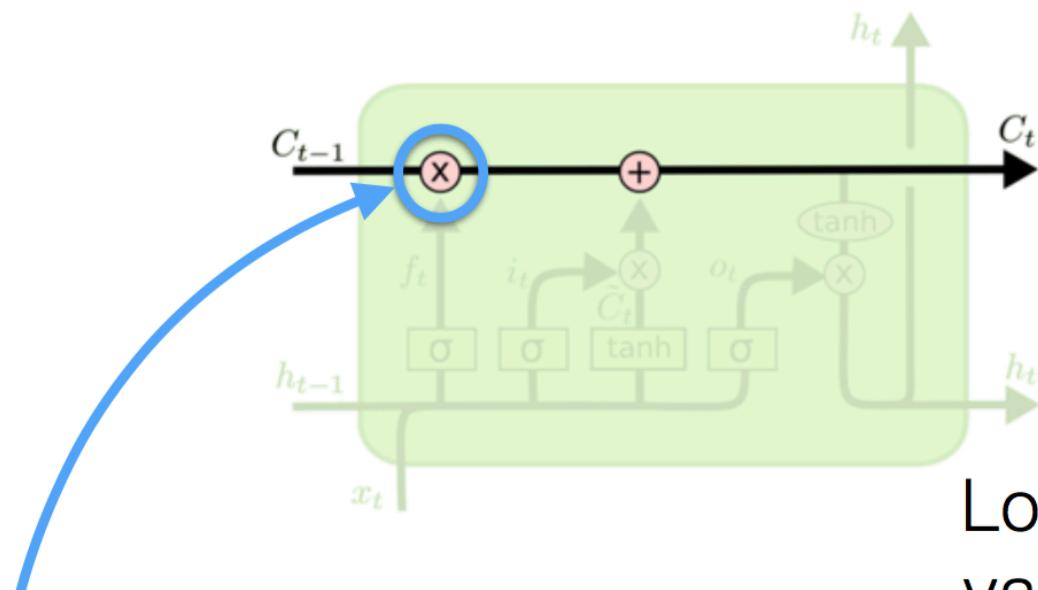


A Gate:



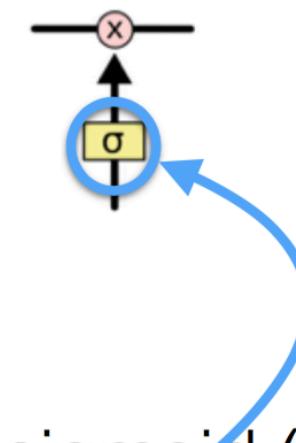
Logistic sigmoid (returns values between 0 and 1)

“The cell”



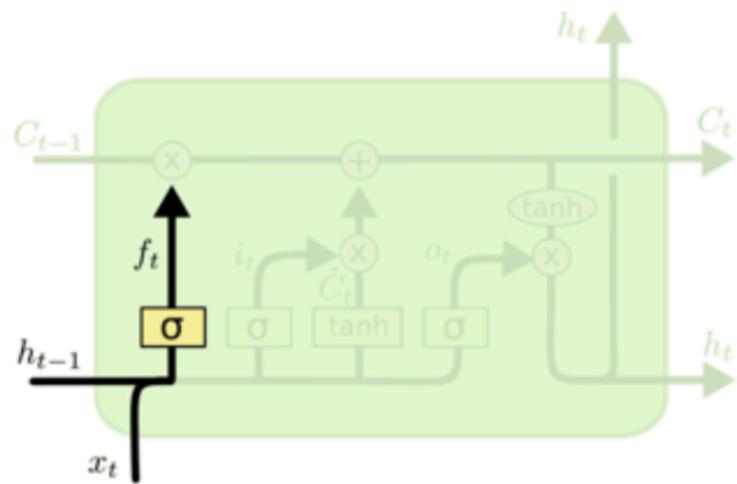
How much of the old  
state should we now forget?

A Gate:



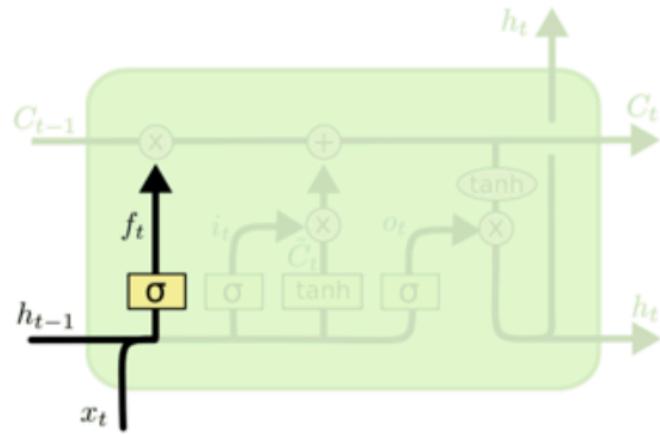
Logistic sigmoid (returns  
values between 0 and 1)

## “Forget gate”



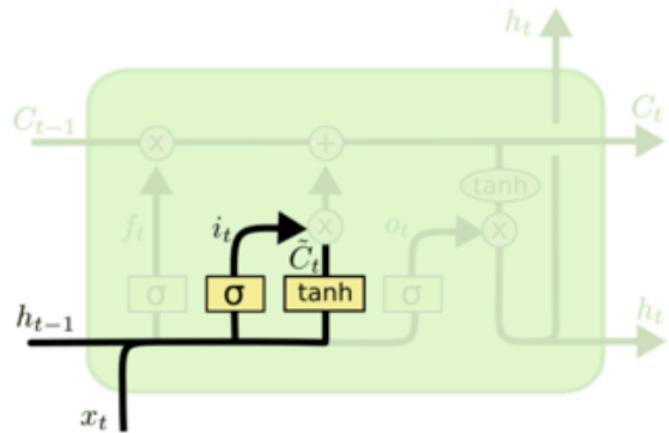
$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

# “Forget gate”



$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

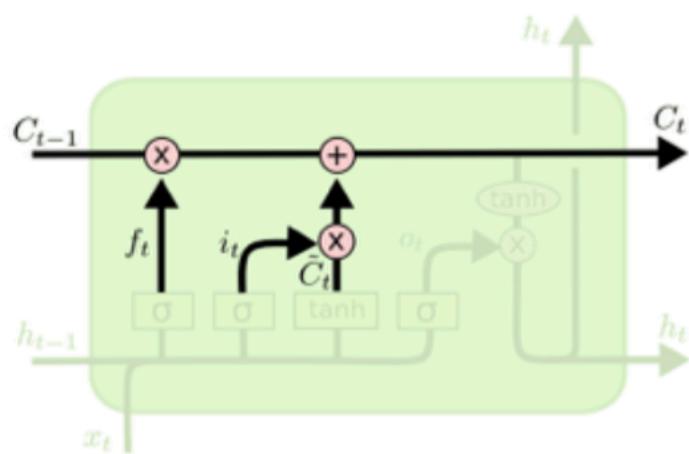
# “Input gate”



$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

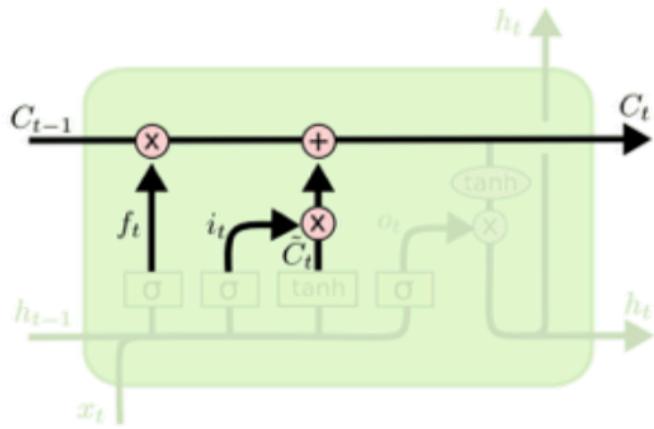
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

# Cell update



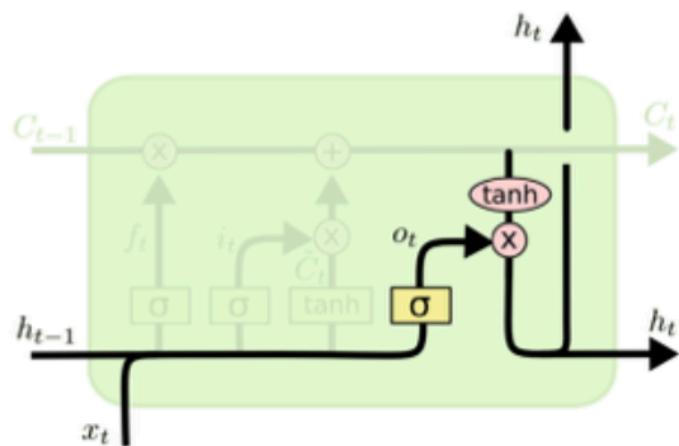
$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

# Cell update



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

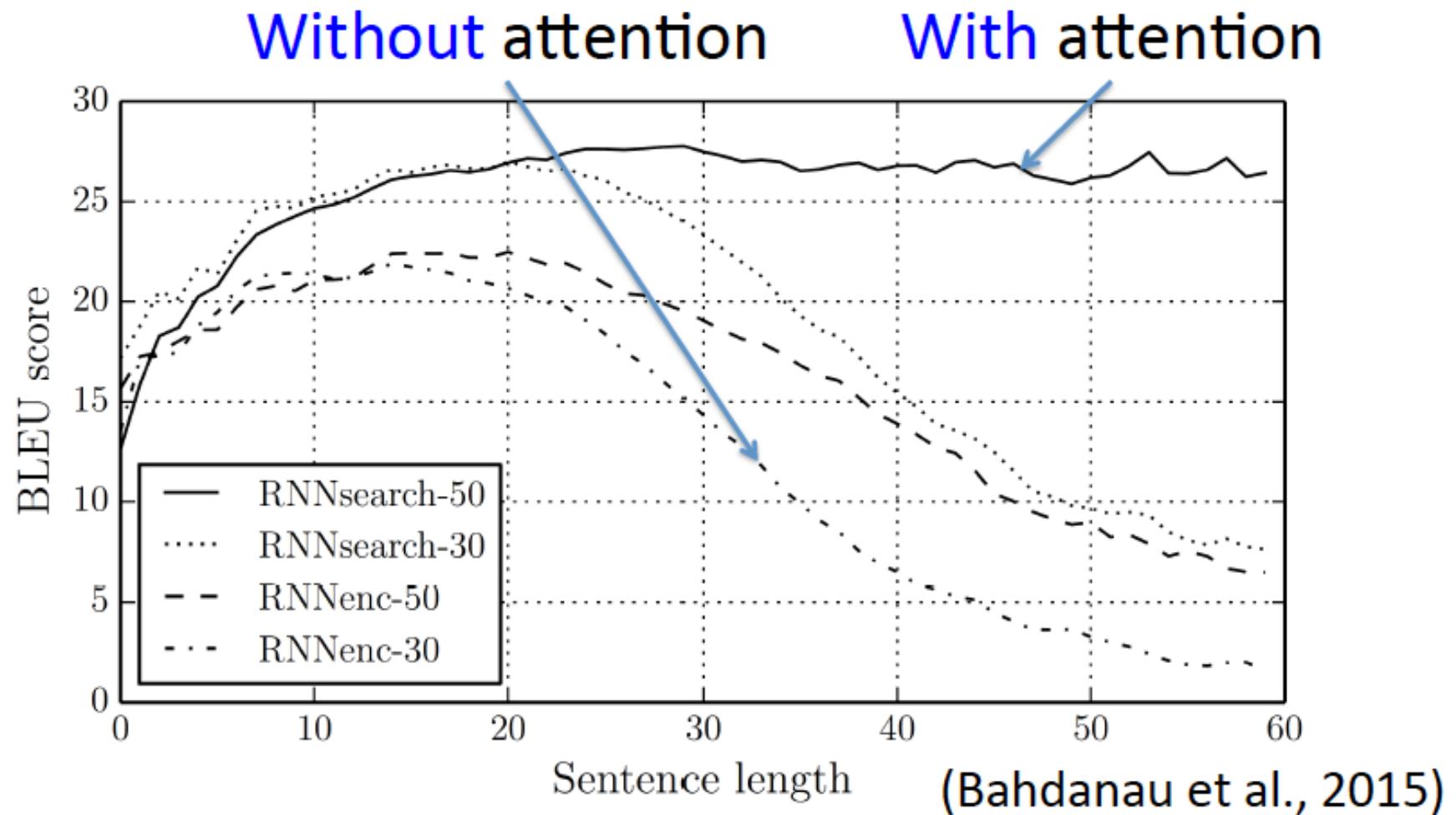
“Output gate” and output



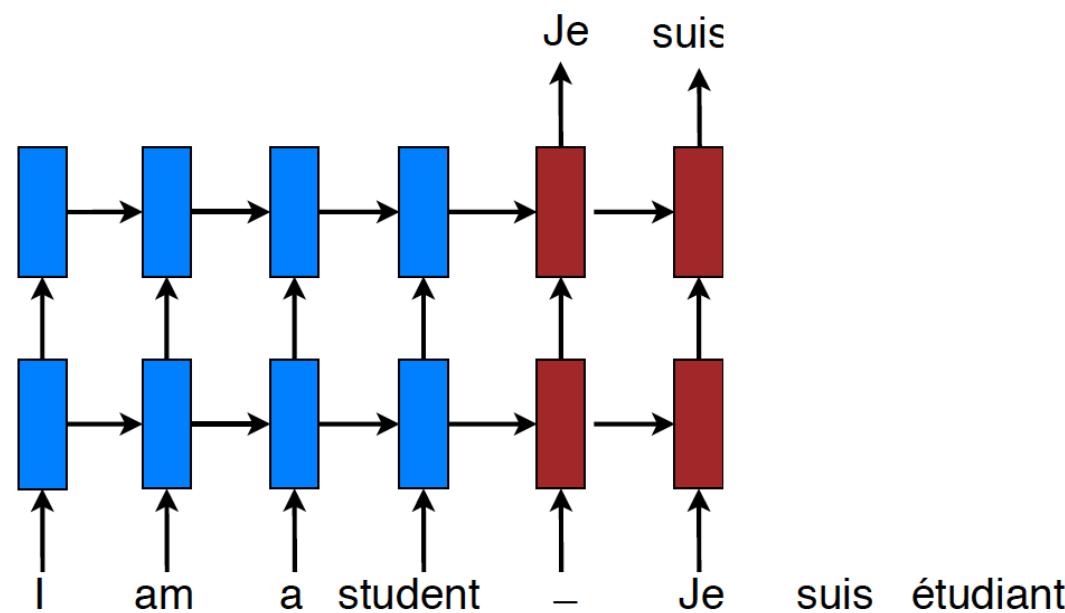
$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh (C_t)$$

# Sentence Length Problem

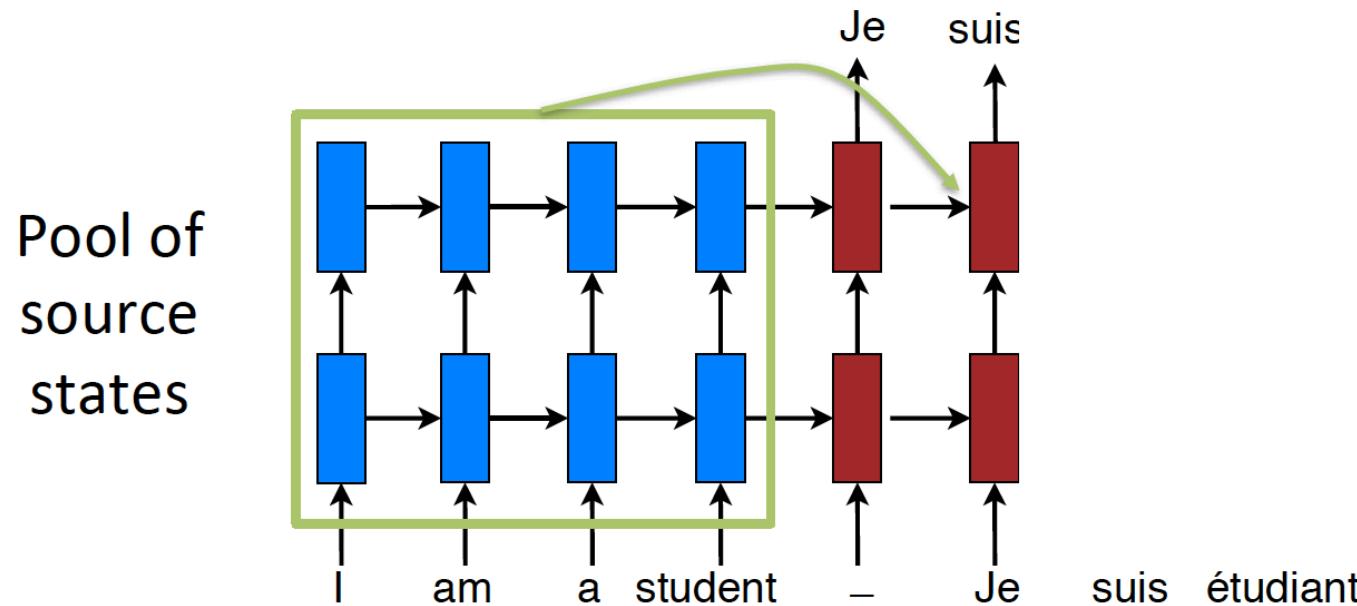


# Why?



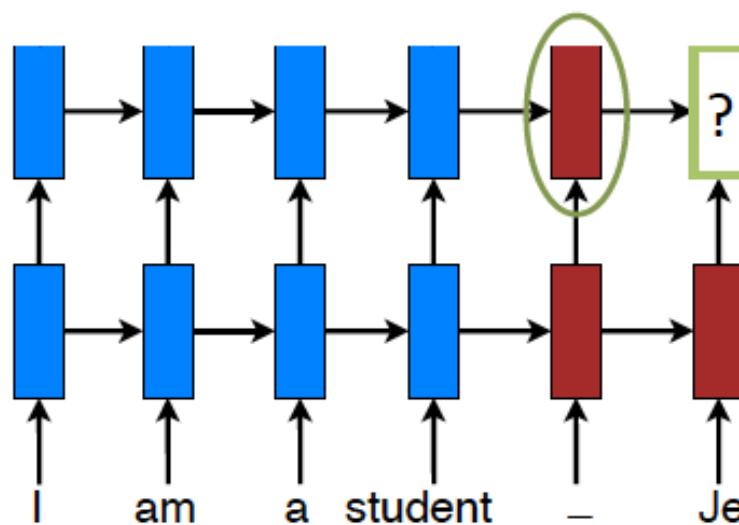
- A **fixed-dimensional** source vector.

# Attention Mechanism



- **Solution:** random access memory
  - Retrieve as needed.

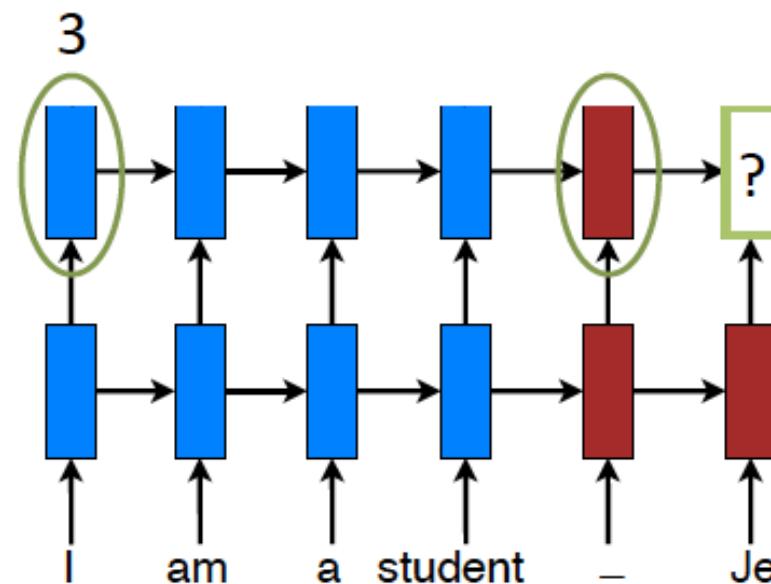
# Attention Mechanism – *Scoring*



- Compare target and source hidden states.

# Attention Mechanism – *Scoring*

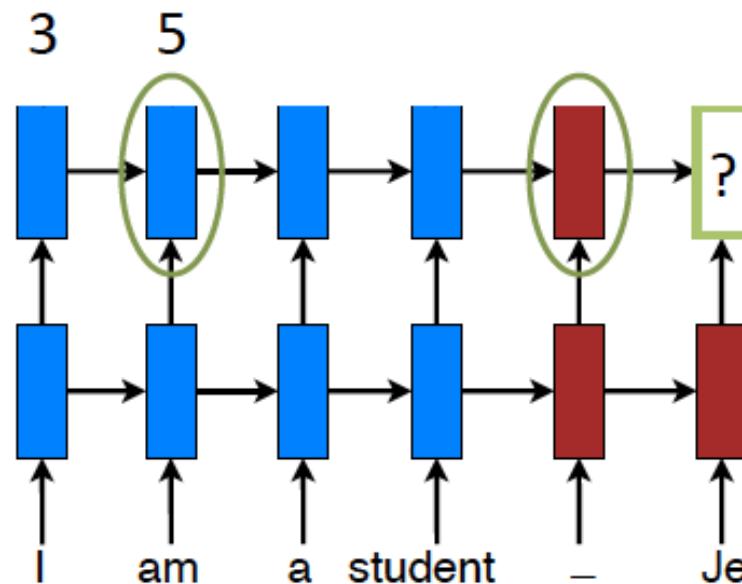
$$\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s)$$



- Compare target and source hidden states.

# Attention Mechanism – *Scoring*

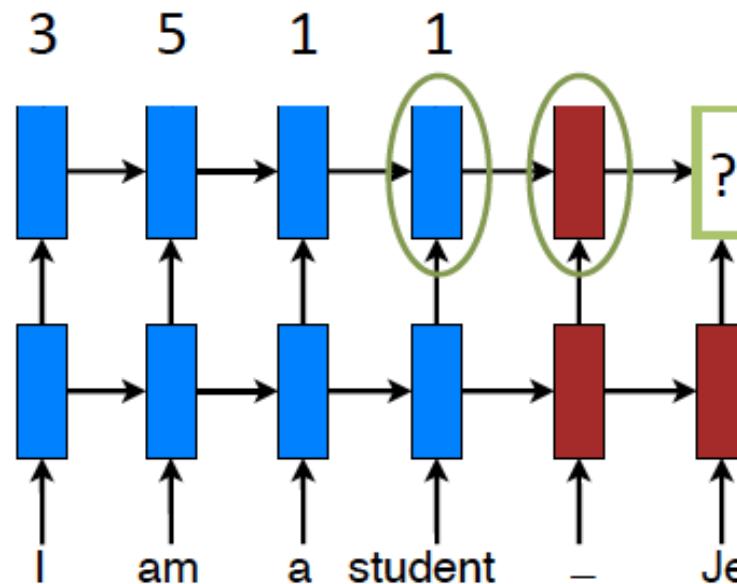
$$\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s)$$



- Compare target and source hidden states.

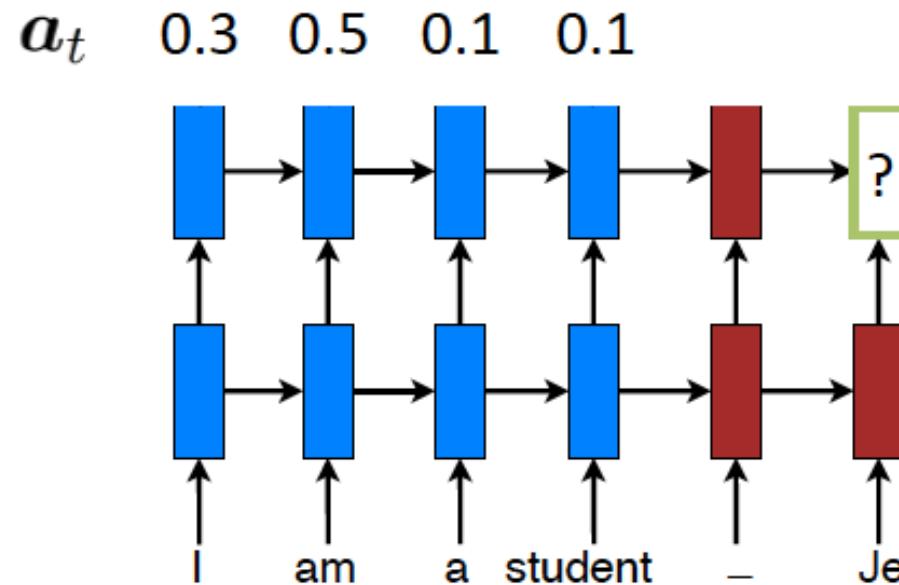
# Attention Mechanism – *Scoring*

$$\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s)$$



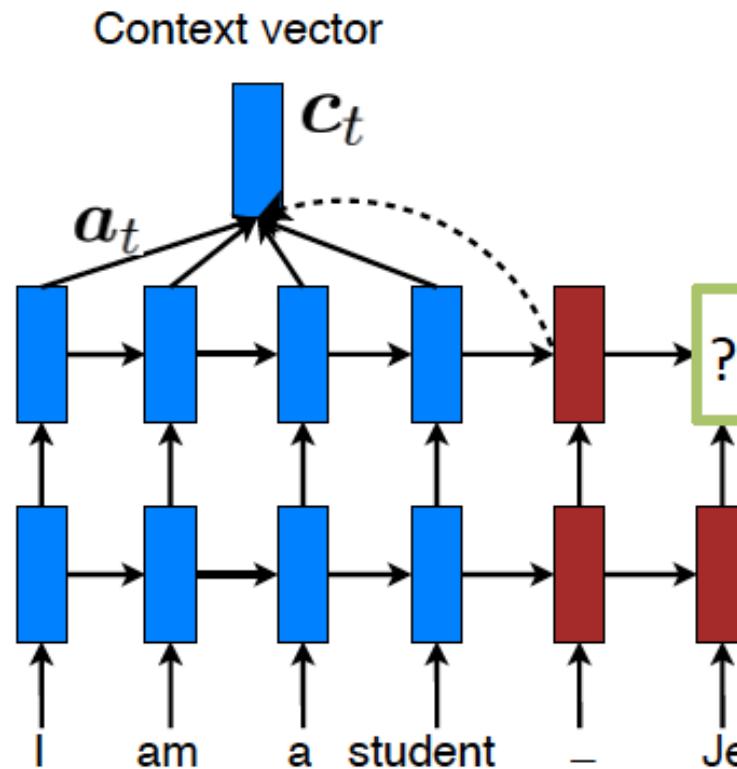
- Compare target and source hidden states.

# Attention Mechanism – *Normalization*



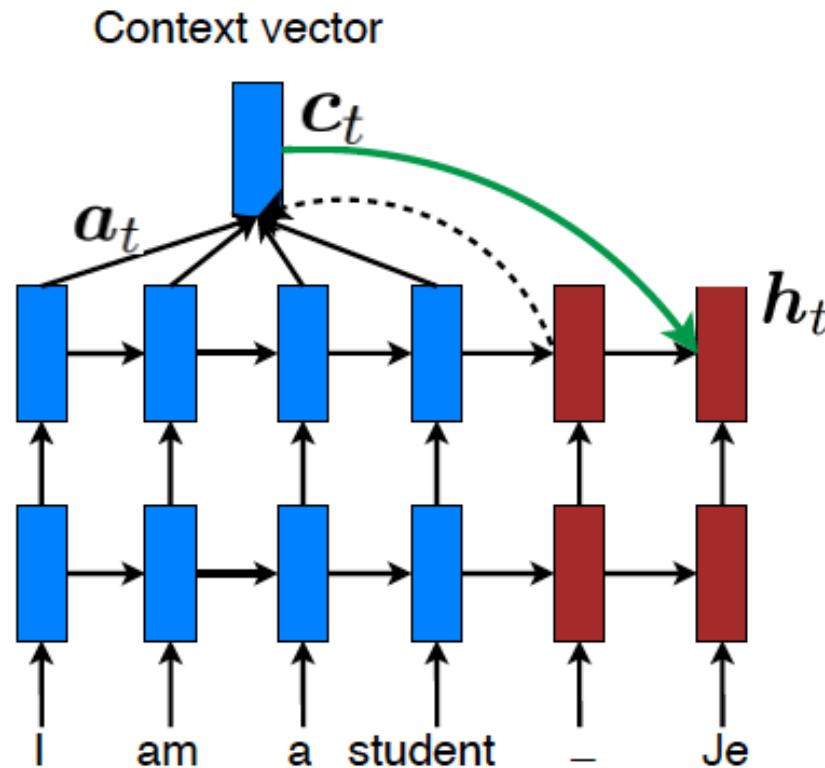
- Convert into alignment weights.

# Attention Mechanism – *Context vector*



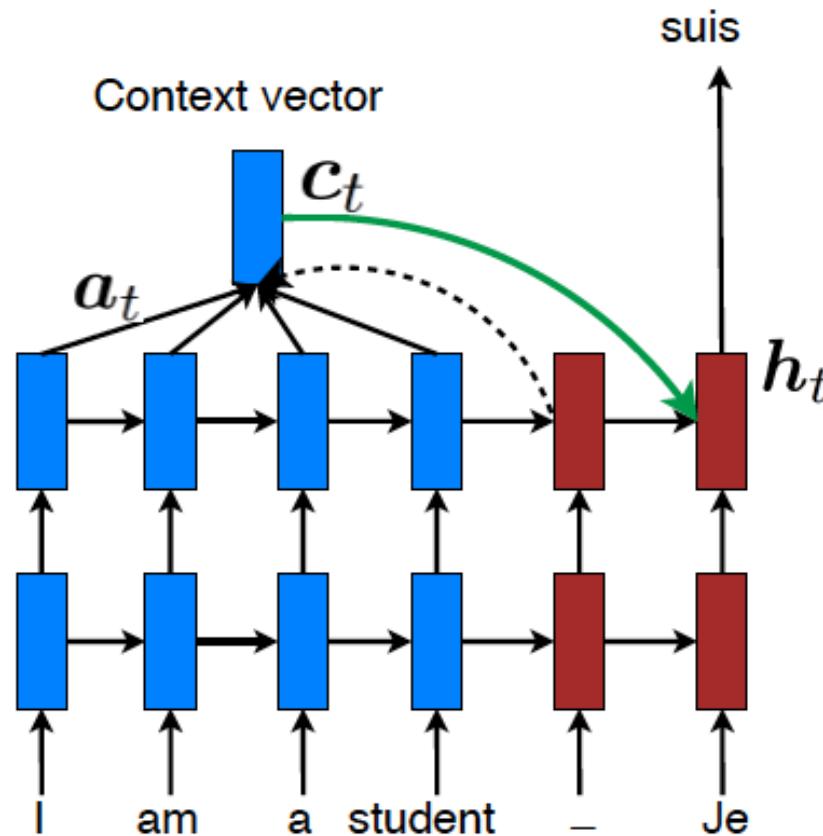
- Build **context** vector: weighted average.

# Attention Mechanism – *Hidden state*



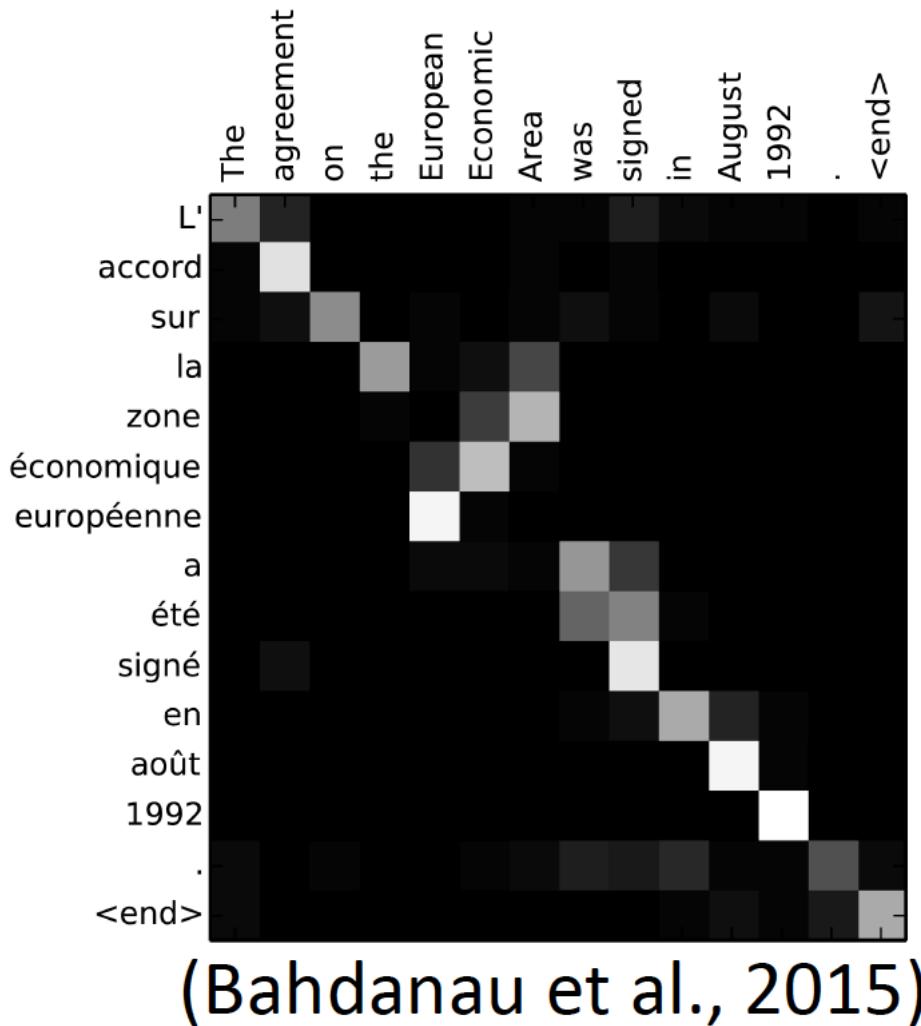
- Compute the next hidden state.

# Attention Mechanism – *Predict*



- Predict the next word.

# Alignments as a by-product



# Acknowledgment

- Chris Dyer (CMU)
- Chris Callison-Burch (JHU)
- Adam Lopez (Edinburgh)
- Philip Koehn
- Thang Luong (Stanford)
- ... many others