

# Lecture Overview

- Machine Translation is one of the big tasks in NLP. Tons of intellectual and technical progress has been driven by the desire to build better MT systems.
- We will focus on “the old way” of doing things today, i.e., Statistical Machine Translation (SMT). Later in the semester, we will talk about neural MT.
- A few key ideas from SMT that are worth knowing:
  - Noisy channel, IBM Model 1, EM algorithm. These are some of the ideas that pushed NLP from rule-based into ML-based models. They are also the kinds of ideas that will come up for other problems (not just NLP).
  - Language models. These are necessary for SMT, and so SMT led to early work on language modeling. Today, language models are all anyone can talk about. ;) We’ll be covering LMs for the next many lectures. Today we will just cover decoding algorithms.

# MT is hard!

- Morphology
  - tuntussuqatarniksaitengqiggtuq (Yupik language, from Alaska)
  - tuntu -ssur -qatar -ni -ksaite -ngqiggte -uq
  - reindeer -hunt -FUTURE -say -NEG -again -3SG.IND
  - "He had not yet said again that he was going to hunt reindeer."

# MT is hard!

- Syntax
  - SVO: German, English, French, Mandarin
  - SOV: Hindi, Japanese
  - VSO: Irish, Arabic, Biblical Hebrew
- Argument structure and marking
  - Possession:
    - The man's house
    - Az ember haza: the man house-his (Hungarian)
  - Motion, manner
    - The bottle floated out
    - La botella salió flotando: The bottle exited floating (Spanish)

# Noisy Channel Model for MT

$$P(tgt | src) = \frac{P(src | tgt)P(tgt)}{P(src)}$$

# Noisy Channel Model for MT

$$P(tgt | src) \propto P(src | tgt)P(tgt)$$

# Noisy Channel Model for MT

$$P(tgt | src) \propto P(src | tgt)P(tgt)$$

Translation Model

# Noisy Channel Model for MT

$$P(tgt | src) \propto P(src | tgt)P(tgt)$$

Translation Model

Intuition: source language is  
a "corrupted" version of the  
target language

# Noisy Channel Model for MT

$$P(tgt | src) \propto P(src | tgt)P(tgt)$$

Language Model



# Noisy Channel Model for MT

$$P(tgt | src) \propto P(src | tgt)P(tgt)$$

Language Model

(We will talk more about  
this in the coming lectures)

# Full Phrase-Based MT System

$$\operatorname{argmax}_{tgt \in TGT} P(src | tgt) P(tgt)$$

# Full Phrase-Based MT System

$$\operatorname{argmax}_{tgt \in TGT} P(src | tgt) P(tgt)$$

## 1. Translation Model

# Full Phrase-Based MT System

$$\operatorname{argmax}_{tgt \in TGT} P(src | tgt) P(tgt)$$

1. Translation Model

2. Language Model

# Full Phrase-Based MT System

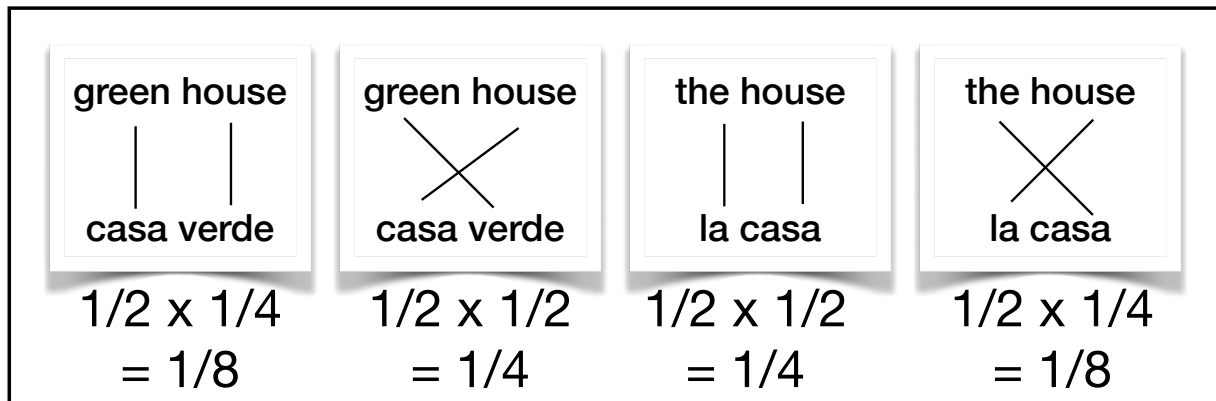
$$\operatorname{argmax}_{tgt \in TGT} P(src | tgt) P(tgt)$$

1. Translation Model
2. Language Model
3. Decoder

# This Lecture: Alignment (Produces a Translation Model)



a



t

		target		
		green	house	the
source	casa	1/2	1/2	1/2
	la	0	1/4	1/2
	verde	1/2	1/4	0

# Lectures 6–8: Language Models

<s> can you tell me about any good cantonese restaurants close by </s>  
    <s> mid priced thai food is what i'm looking for </s>  
        <s> tell me about chez panisse </s>  
<s> can you give me a listing of the kinds of food that are available </s>  
    <s> i'm looking for a good place to eat breakfast </s>  
        <s> when is caffe venezia open during the day </s>

$$P(w_0 \dots w_n) \approx P(w_0 | \text{<s>}) \times P(w_1 | w_0) \times P(w_2 | w_1) \times \dots \times P(w_n | w_{n-1})$$

$$P(\text{tell me about caffe venezia}) = P(\text{tell} | \text{<s>}) \times P(\text{me} | \text{tell}) \times P(\text{about} | \text{me}) \times \\ P(\text{caffe} | \text{about}) \times P(\text{venezia} | \text{caffe}) \times P(\text{</s>} | \text{venezia})$$

# Full Phrase-Based MT System

$$\operatorname{argmax}_{tgt \in TGT} P(src | tgt) P(tgt)$$

1. Translation Model
2. Language Model
3. Decoder



# Word Alignment

Goal: Build a phrase table

t		target		
		green	house	the
source	casa	1/2	1/2	1/2
	la	0	1/4	1/2
	verde	1/2	1/4	0

# Word Alignment

## Bitexts a.k.a. Bilingual Parallel Corpora

- Bitext: A corpus that contains translated pairs of texts
  - Can be aligned coarsely (e.g., document-level) or finely (sentence level)
  - Word-level alignment is rare

# Word Alignment

## Bitexts a.k.a. Bilingual Parallel Corpora

- Some exist naturally. E.g., EU translates all their documents (<https://www.statmt.org/europarl/>), authors hire translators to translate literature

Parallel Corpus (L1-L2)	Sentences	L1 Words	English Words
Bulgarian-English	406,934	-	9,886,291
Czech-English	646,605	12,999,455	15,625,264
Danish-English	1,968,800	44,654,417	48,574,988
German-English	1,920,209	44,548,491	47,818,827
Greek-English	1,235,976	-	31,929,703
Spanish-English	1,965,734	51,575,748	49,093,806
Estonian-English	651,746	11,214,221	15,685,733
Finnish-English	1,924,942	32,266,343	47,460,063
French-English	2,007,723	51,388,643	50,196,035
Hungarian-English	624,934	12,420,276	15,096,358
Italian-English	1,909,115	47,402,927	49,666,692
Lithuanian-English	635,146	11,294,690	15,341,983
Latvian-English	637,599	11,928,716	15,411,980
Dutch-English	1,997,775	50,602,994	49,469,373
Polish-English	632,565	12,815,544	15,268,824
Portuguese-English	1,960,407	49,147,826	49,216,896
Romanian-English	399,375	9,628,010	9,710,331

# Word Bitexts

- Some all their documents [www.statmt.org/eu](http://www.statmt.org/eu) translators to translate
- Can get creative if we are okay with “weakly aligned” —e.g., pairs of new articles or Wikipedia documents on the same topic

## Natural language processing

From Wikipedia, the free encyclopedia

*This article is about natural language processing done by computers. For the natural language processing done by the human brain, see [Language processing in the brain](#).*

**Natural language processing (NLP)** is a subfield of [linguistics](#), [computer science](#), and [artificial intelligence](#) concerned with the interactions between computers and human language, in particular how to program computers to process and analyze large amounts of [natural language](#) data. The goal is a computer capable of “understanding” the contents of documents, including the [contextual](#) nuances of the language within them. The technology can then accurately extract information and insights contained in the documents as well as categorize and organize the documents themselves.

Challenges in natural language processing frequently

### Contents [\[hide\]](#)

#### 1 History

- 1.1 Symbolic NLP (1950s – early 1990s)
- 1.2 Statistical NLP (1990s–2010s)
- 1.3 Neural NLP (present)

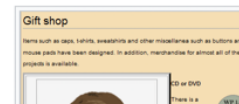
## Procesamiento de lenguajes naturales

(Redirigido desde «NLP»)

El **procesamiento de lenguaje natural**,<sup>1 2</sup> abreviado **PLN**<sup>3 4</sup> —en inglés, *natural language processing*, NLP— es un campo de las [ciencias de la computación](#), de la [inteligencia artificial](#) y de la [lingüística](#) que estudia las interacciones entre las computadoras y el lenguaje humano. Se ocupa de la formulación e investigación de mecanismos eficaces computacionalmente para la comunicación entre personas y máquinas por medio del [lenguaje natural](#), es decir, de las [lenguas del mundo](#). No trata de la comunicación por medio de lenguas naturales de una forma abstracta, sino de diseñar mecanismos para comunicarse que sean eficaces computacionalmente —que se puedan realizar por medio de programas que ejecuten o simulen la comunicación—. Los modelos aplicados se enfocan no solo a la comprensión del lenguaje de por sí, sino a aspectos generales cognitivos humanos y a la organización de la memoria. El lenguaje natural sirve solo de medio para estudiar estos fenómenos. Hasta la década de 1980, la mayoría de los sistemas de PLN se basaban en un complejo conjunto de reglas diseñadas a mano. A partir de finales de 1980, sin embargo, hubo una revolución en PLN con la introducción de [algoritmos de aprendizaje automático](#) para el procesamiento del lenguaje.<sup>5 6</sup>

### Índice [\[ocultar\]](#)

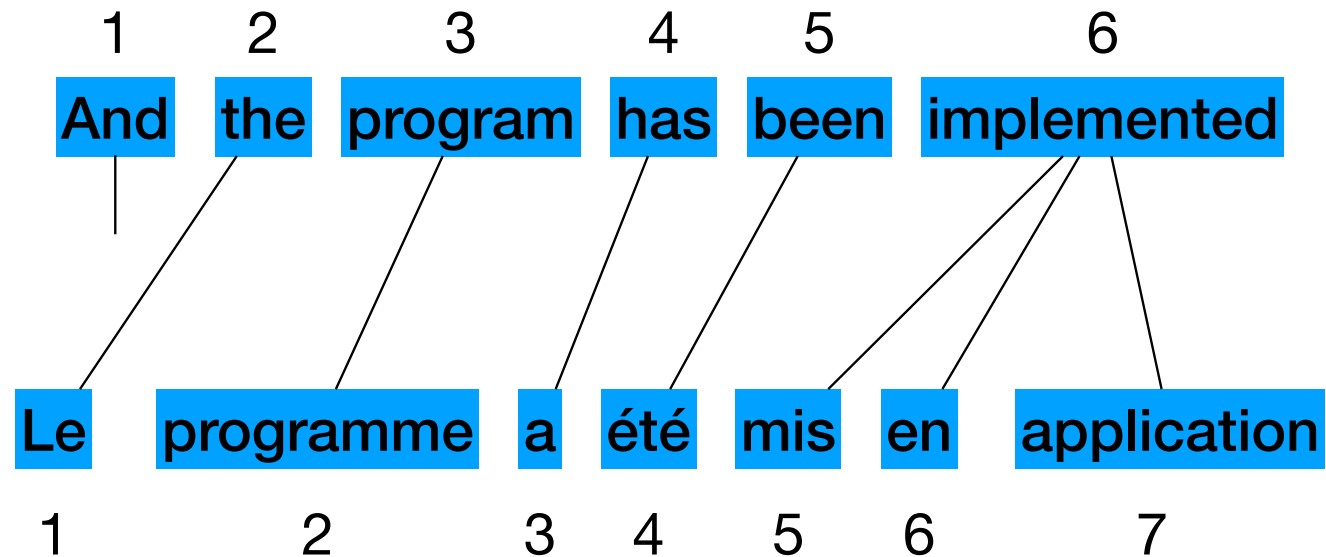
- 1 Historia
- 2 Dificultades en el procesamiento de lenguaje natural
  - 2.1 Ambigüedad



Estonian-English	651,740	11,214,221	15,003,733
Finnish-English	1,924,942	32,266,343	47,460,063
French-English	2,007,723	51,388,643	50,196,035
Hungarian-English	624,934	12,420,276	15,096,358
Italian-English	1,909,115	47,402,927	49,666,692
Lithuanian-English	635,146	11,294,690	15,341,983
Latvian-English	637,599	11,928,716	15,411,980
Dutch-English	1,997,775	50,602,994	49,469,373
Polish-English	632,565	12,815,544	15,268,824
Portuguese-English	1,960,407	49,147,826	49,216,896
Romanian-English	399,375	9,628,010	9,710,331

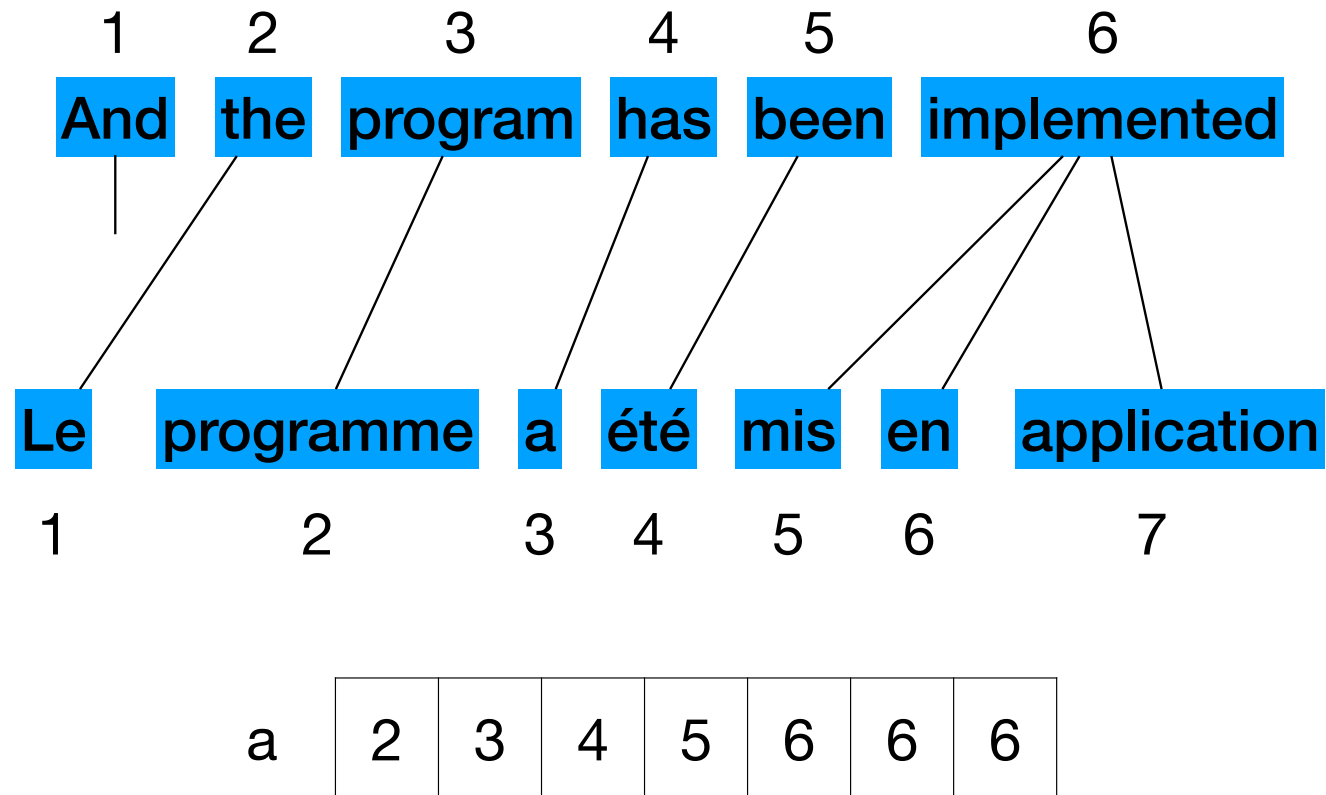
# Word Alignment

## IBM Model 1



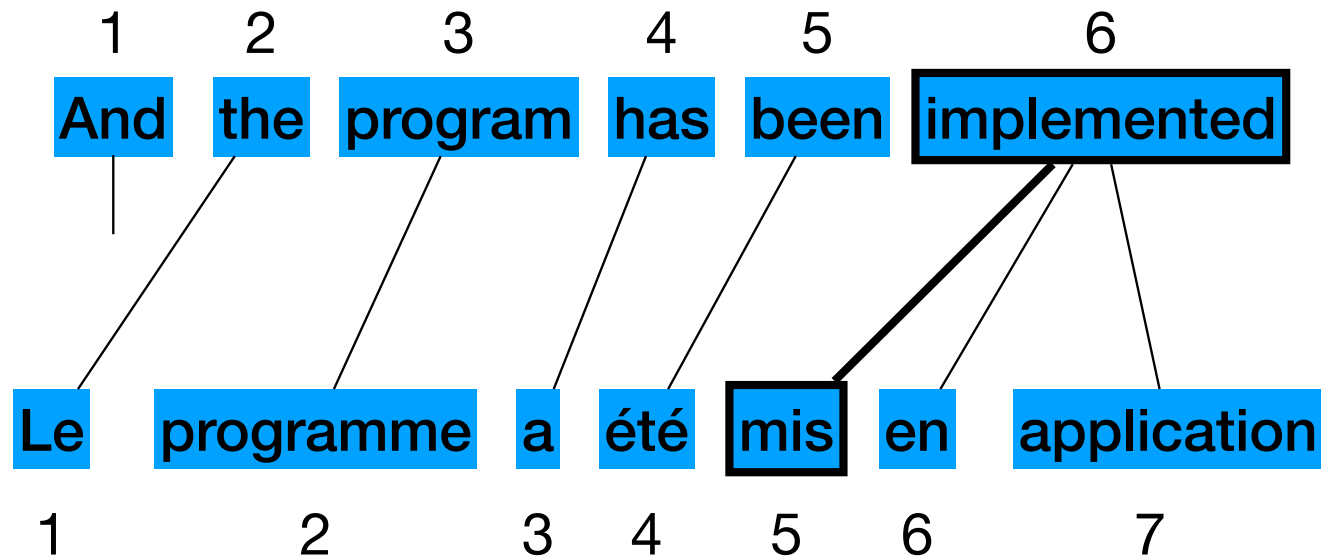
# Word Alignment

## IBM Model 1



# Word Alignment

## IBM Model 1

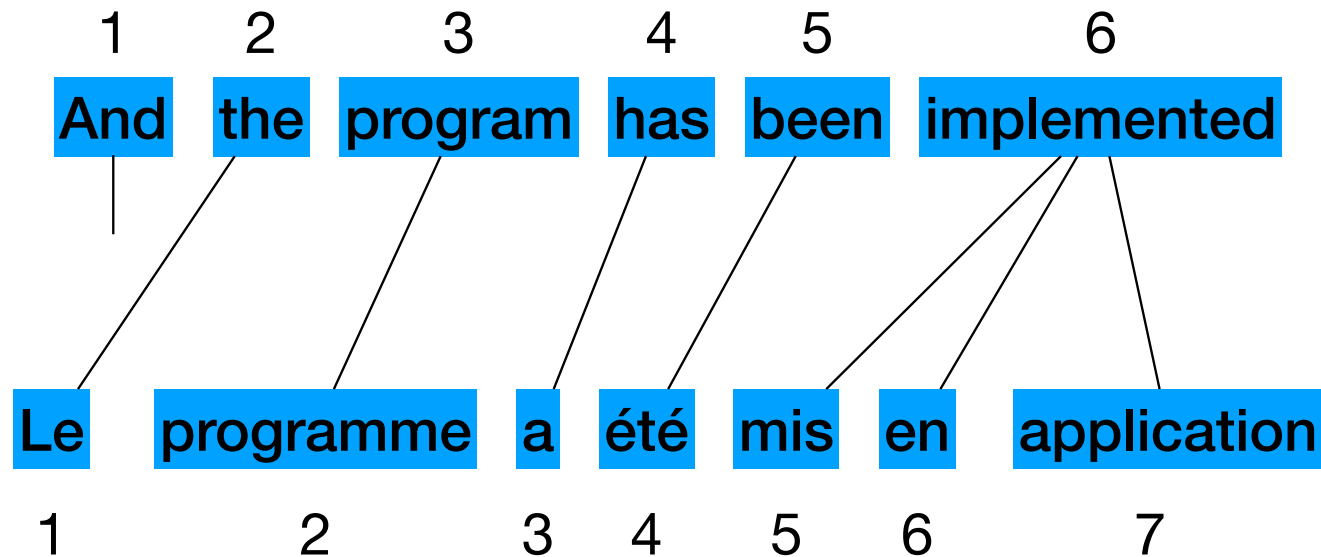


a	2	3	4	5	6	6	6
---	---	---	---	---	---	---	---

$$a[5] = 6$$

# Word Alignment

## IBM Model 1



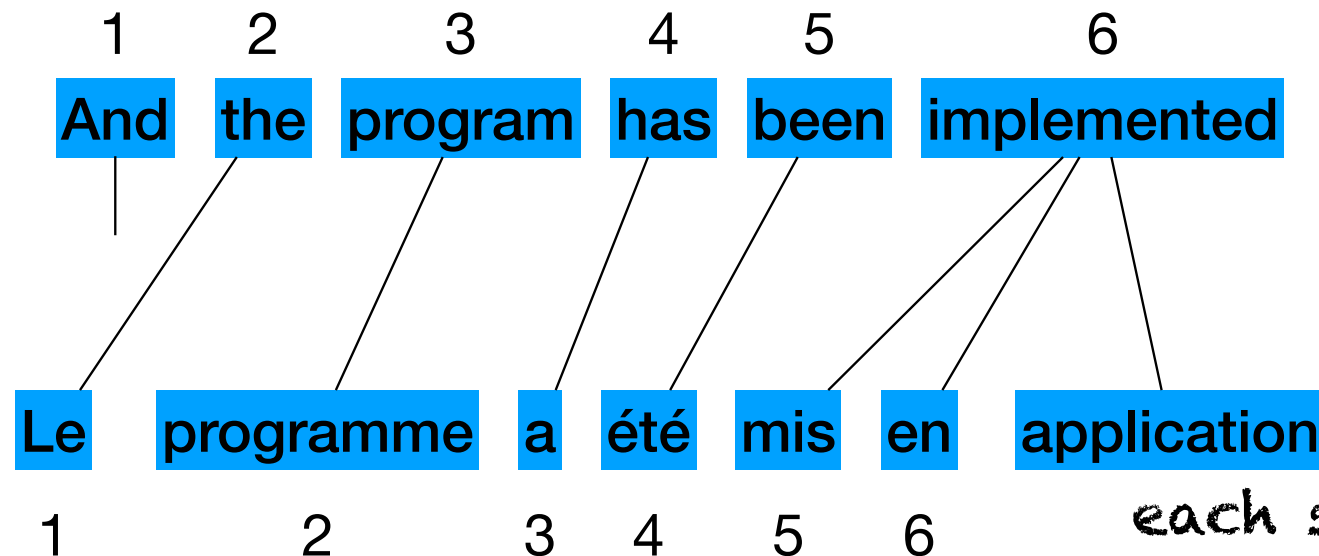
a	2	3	4	5	6	6	6
---	---	---	---	---	---	---	---

each source word  
aligns to exactly  
one target word



# Word Alignment

## IBM Model 1



a	2	3	4	5	6	6	6
---	---	---	---	---	---	---	---

each source word  
aligns to exactly  
one target word  
(later models  
relax this)

# Word Alignment

## IBM Model 1

- First (simplest) automatic unsupervised alignment model from IBM
- Generative Story:
  - Choose length  $L$  for src sentence
  - Choose an alignment  $A = a_1 \dots a_L$
  - Then, generate target position  $t_i$  by translating whatever source phrase is aligned to position  $i$  in the target
- To build the phrase table:  $\operatorname{argmax}_{aj} (s_j \mid t_{aj})$

# Word Alignment

## IBM Model 1

- Dilemma:
  - We want alignments so that we can figure out the probabilities of phrase translations, e.g.,  $P(s_i|t_j)$  for all  $s, t, i, j$
  - To estimate those alignments, we need phrase translation probabilities
  - 😞
  - No fear! EM is here!

# Topics

- Follow ups
  - SVD Revisited
  - Generative Stories: Intuition
- Machine Translation
- Noisy Channel Models for SMT
  - **Translation Model (Word Alignment)**
    - IBM Alignment Models
    - **EM Algorithm**
  - Language Model
    - Decoding

# Expectation Maximization (EM) Algorithm

- E step: Estimate the likelihood of the observed data given the current parameters
- M step: Recompute the parameters so as to maximize the likelihood of the observed data

# Expectation Maximization (EM) Algorithm

Goal:

Estimate a and t using

$$\operatorname{argmax}_{a_j} (s_j \mid t_{a_j})$$

green house

the house

casa verde

la casa

		target		
		green	house	the
source	casa	1/3	1/3	1/3
	la	1/3	1/3	1/3
	verde	1/3	1/3	1/3

# Expectation Maximization (EM) Algorithm

Goal:

Estimate a and t using

$$\operatorname{argmax}_{a_j} (s_j \mid t_{a_j})$$

Corpus

green house	casa verde
the house	la casa

		target		
		green	house	the
source	casa	1/3	1/3	1/3
	la	1/3	1/3	1/3
	verde	1/3	1/3	1/3

# Expectation Maximization (EM) Algorithm

Goal:

Estimate a and t using

$$\operatorname{argmax}_{a_j} (s_j \mid t_{a_j})$$

Parameters

green house

the house

casa verde

la casa

t

		target		
		green	house	the
source	casa	1/3	1/3	1/3
	la	1/3	1/3	1/3
	verde	1/3	1/3	1/3



# Expectation Maximization (EM) Algorithm

Goal:

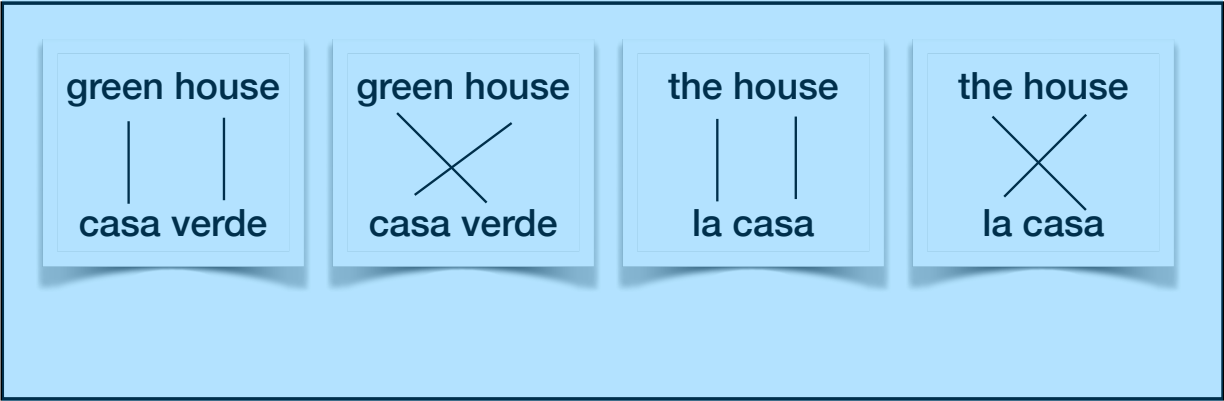
Estimate  $a$  and  $t$  using

$$\operatorname{argmax}_{a_j} (s_j \mid t_{aj})$$

All possible alignments



$a$



$t$

		target		
		green	house	the
source	casa	1/3	1/3	1/3
	la	1/3	1/3	1/3
	verde	1/3	1/3	1/3

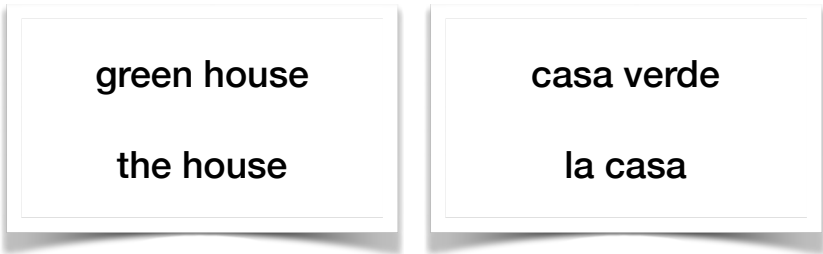
# Expectation Maximization (EM) Algorithm

Goal:

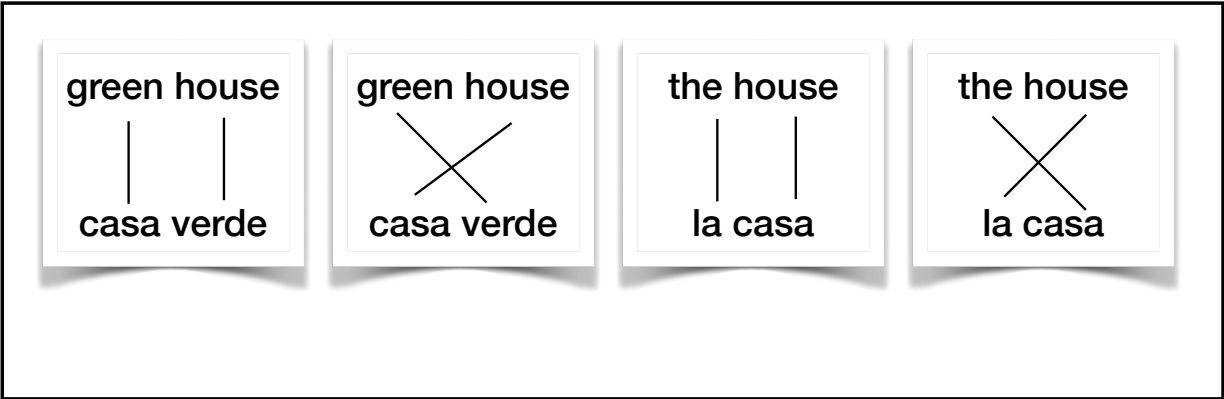
Estimate  $a$  and  $t$  using

$$\operatorname{argmax}_{a_j} (s_j \mid t_{aj})$$

E Step: Compute Likelihood of corpus given current parameters



a



		target		
		green	house	the
source	casa	1/3	1/3	1/3
	la	1/3	1/3	1/3
	verde	1/3	1/3	1/3

# Expectation Maximization (EM) Algorithm

Goal:

Estimate  $a$  and  $t$  using

$$\operatorname{argmax}_{a_j} (s_j \mid t_{aj})$$

E Step: Compute Likelihood of corpus given current parameters

green house

the house

casa verde

la casa

a

green house

  
casa verde

$$\frac{1}{3} \times \frac{1}{3} = \frac{1}{9}$$

green house

  
casa verde

the house

  
la casa

the house

  
la casa

t

target

green

house

the

source

casa

$\frac{1}{3}$

$\frac{1}{3}$

$\frac{1}{3}$

la

$\frac{1}{3}$

$\frac{1}{3}$

$\frac{1}{3}$

verde

$\frac{1}{3}$

$\frac{1}{3}$

$\frac{1}{3}$

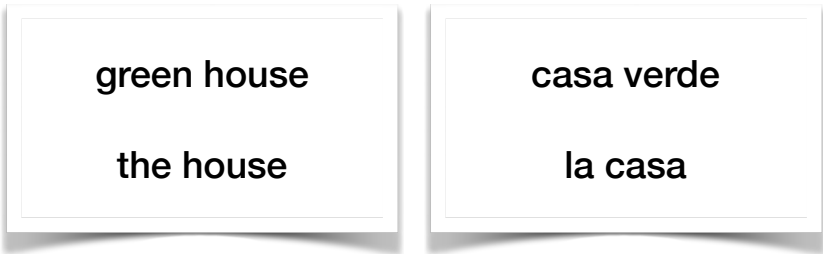
# Expectation Maximization (EM) Algorithm

Goal:

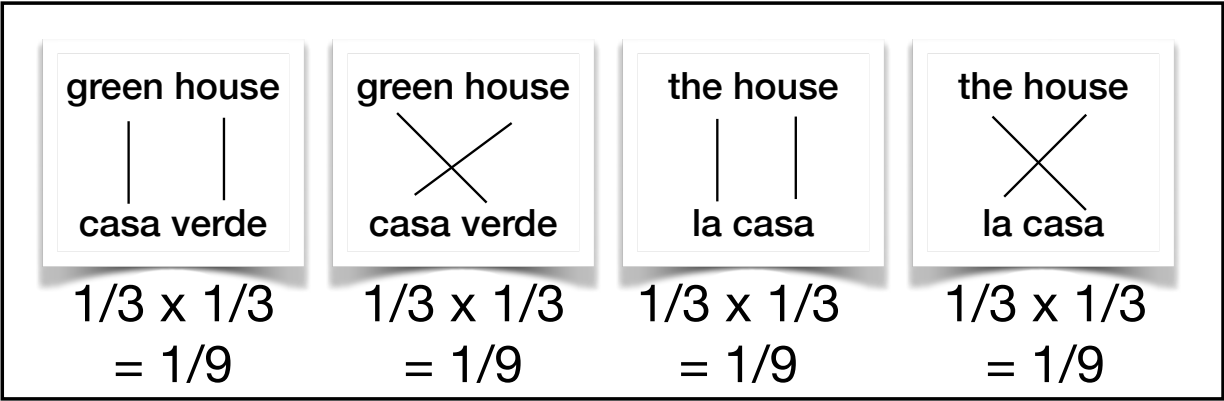
Estimate  $a$  and  $t$  using

$$\operatorname{argmax}_{a_j} (s_j \mid t_{aj})$$

E Step: Compute Likelihood of corpus given current parameters



a



		target		
		green	house	the
source	casa	1/3	1/3	1/3
	la	1/3	1/3	1/3
	verde	1/3	1/3	1/3

# Expectation Maximization (EM) Algorithm

Goal:

Estimate  $a$  and  $t$  using

$$\operatorname{argmax}_{a_j} (s_j \mid t_{aj})$$

E Step: Compute Likelihood of corpus given current parameters

$$P(a \mid s, t) = \frac{P(a, s \mid t)}{\sum_a P(a, s \mid t)}$$

Normalize to reweight likelihood of alignments

green house	casa verde
the house	la casa

$a$

<div>green house</div> <div> <div></div> <div></div> </div> <div>casa verde</div> <div>=1/2</div>	<div>green house</div> <div> <div></div> <div></div> </div> <div>casa verde</div> <div>=1/2</div>	<div>the house</div> <div> <div></div> <div></div> </div> <div>la casa</div> <div>=1/2</div>	<div>the house</div> <div> <div></div> <div></div> </div> <div>la casa</div> <div>=1/2</div>
---	---	--	--

$t$

		target		
		green	house	the
source	casa	1/3	1/3	1/3
	la	1/3	1/3	1/3
	verde	1/3	1/3	1/3

# Expectation Maximization (EM) Algorithm

Goal:

Estimate  $a$  and  $t$  using

$$\operatorname{argmax}_{a_j} (s_j | t_{aj})$$

E Step: Compute Likelihood of corpus given current parameters

Compute expected counts counts by adding fractional counts equal to  $P(a|f,e)$ .

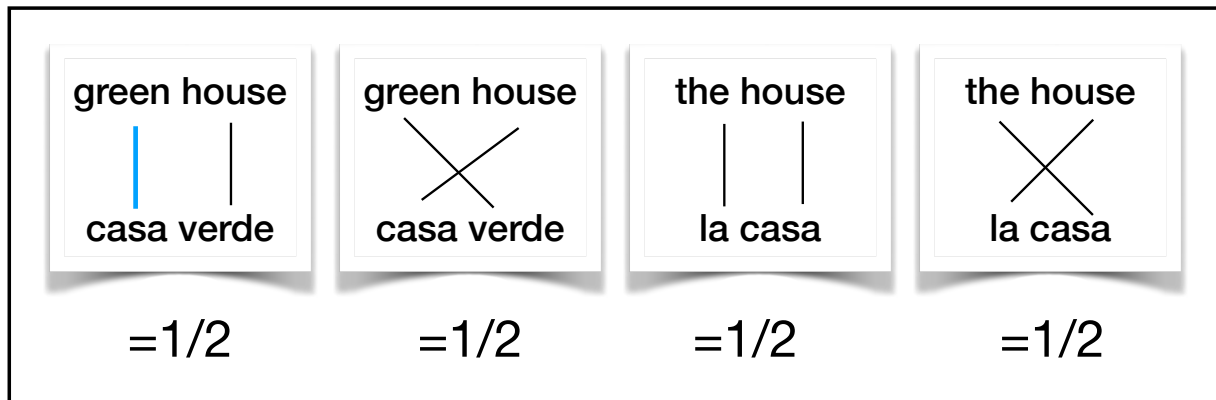
green house

the house

casa verde

la casa

a



t

		target		
		green	house	the
source	casa	1/2	1	1/2
	la	0	1/2	1/2
	verde	1/2	1/2	0

# Expectation Maximization (EM) Algorithm

Goal:

Estimate  $a$  and  $t$  using

$$\operatorname{argmax}_{a_j} (s_j \mid t_{aj})$$

E Step: Compute Likelihood of corpus given current parameters

Compute expected counts by adding fractional counts equal to  $P(a|f,e)$ .

green house	casa verde
the house	la casa

a

<div>green house   casa verde =1/2</div>	<div>green house <del> </del> casa verde =1/2</div>	<div>the house   la casa =1/2</div>	<div>the house <del> </del> la casa =1/2</div>
--	---	---	--

		target		
		green	house	the
source	casa	1/2	1	1/2
	la	0	1/2	1/2
	verde	1/2	1/2	0

# Expectation Maximization (EM) Algorithm

Goal:

Estimate  $a$  and  $t$  using

$$\operatorname{argmax}_{a_j} (s_j \mid t_{aj})$$

*M Step: Compute MLE parameters*

*(here, that just means normalizing!)*

green house	casa verde
the house	la casa

*a*

green house	green house	the house	the house
casa verde	<del>casa verde</del>	la casa	<del>la casa</del>

		target		
		green	house	the
source	casa	1/2	1/2	1/2
	la	0	1/4	1/2
	verde	1/2	1/4	0



# Expectation Maximization (EM) Algorithm

Goal:

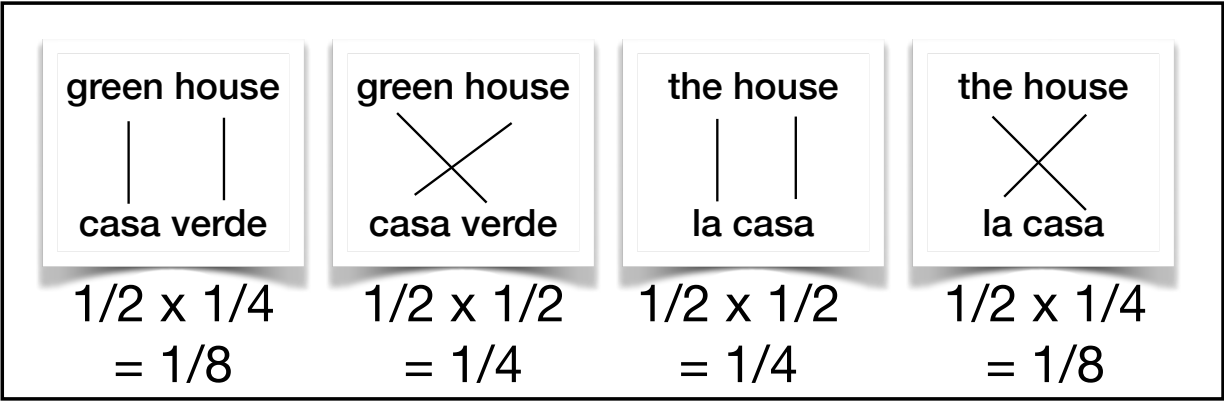
Estimate  $a$  and  $t$  using

$$\operatorname{argmax}_{a_j} (s_j \mid t_{aj})$$

E Step: Compute Likelihood of corpus given current parameters



a



		target		
		green	house	the
source	casa	1/2	1/2	1/2
	la	0	1/4	1/2
	verde	1/2	1/4	0

# Expectation Maximization (EM) Algorithm

Goal:

Estimate  $a$  and  $t$  using

$$\operatorname{argmax}_{a_j} (s_j \mid t_{aj})$$

E Step: Compute Likelihood of corpus given current parameters

green house

the house

casa verde

la casa

a

<div>green house</div> <div> <div></div> <div></div> </div> <div>casa verde</div> <div> <math>1/2 \times 1/4</math>  <math>= 1/8</math> </div>	<div>green house</div> <div> <div></div> <div></div> </div> <div>casa verde</div> <div> <math>1/2 \times 1/2</math>  <math>= 1/4</math> </div>	<div>the house</div> <div> <div></div> <div></div> </div> <div>la casa</div> <div> <math>1/2 \times 1/2</math>  <math>= 1/4</math> </div>	<div>the house</div> <div> <div></div> <div></div> </div> <div>la casa</div> <div> <math>1/2 \times 1/4</math>  <math>= 1/8</math> </div>
--	--	---	---

t

		target		
		green	house	the
source	casa	1/2	1/2	1/2
	la	0	1/4	1/2
	verde	1/2	1/4	0