

Machine Translation (Part I)

CS 585, Fall 2015

Introduction to Natural Language Processing

<http://people.cs.umass.edu/~brenocon/inlp2015/>

Brendan O'Connor

College of Information and Computer Sciences

University of Massachusetts Amherst

[Some slides borrowed from J&M + mt-class.org]

NLP news of the day

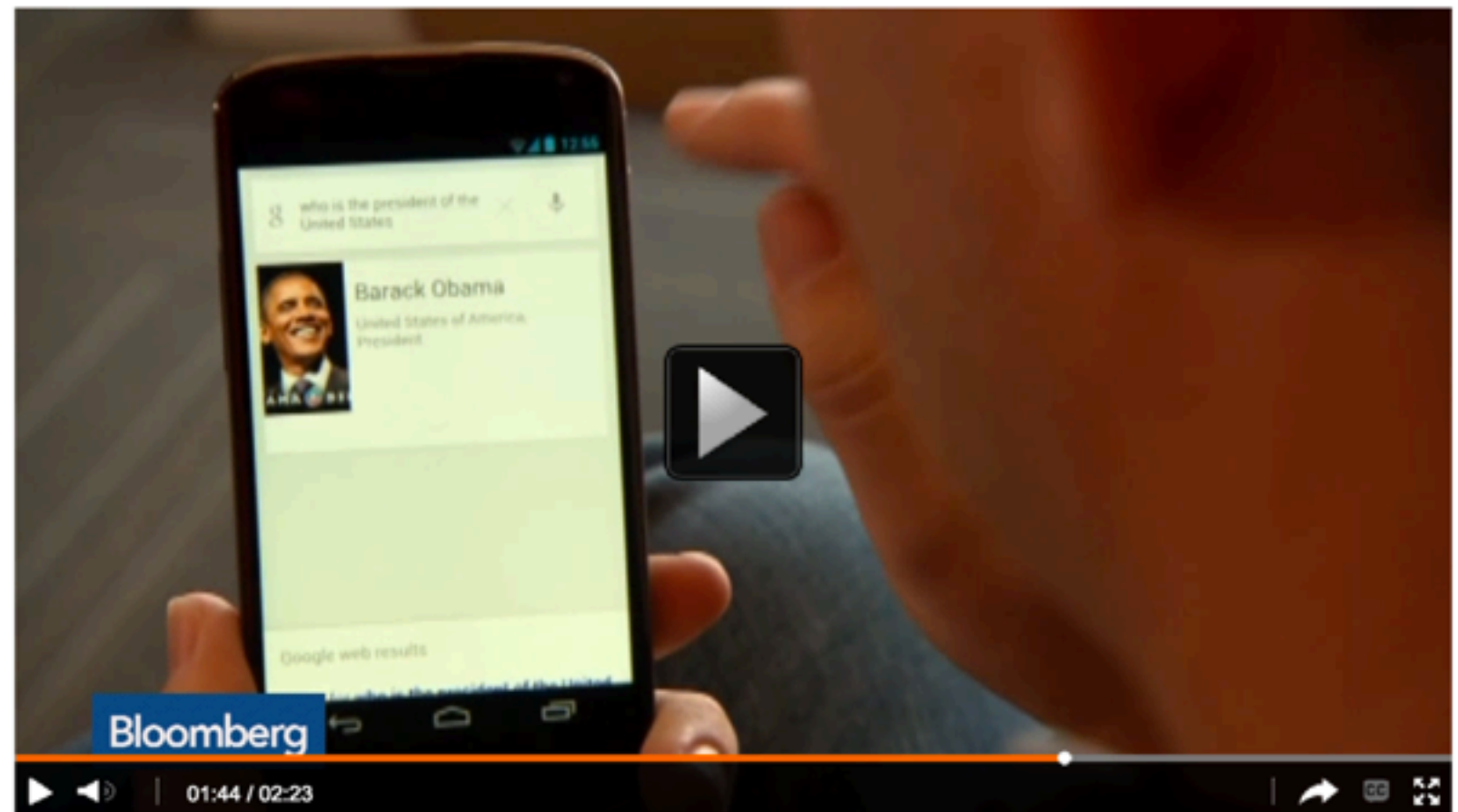
The system helps Mountain View, California-based Google deal with the 15 percent of queries a day it gets which its systems have never seen before, he said. For example, it's adept at dealing with ambiguous queries, like, "What's the title of the consumer at the highest level of a food chain?" And RankBrain's usage of AI means it works differently than the other technologies in the search engine.

"The other signals, they're all based on discoveries and insights that people in information retrieval have had, but there's no learning," Corrado said.

Google Turning Its Lucrative Web Search Over to AI Machines

by Jack Clark

October 26, 2015 – 5:00 AM EDT



Graders at work

- Midterms back on Thursday
- Project feedback: by tomorrow
- HW2 still underway (sorry!)

Machine translation

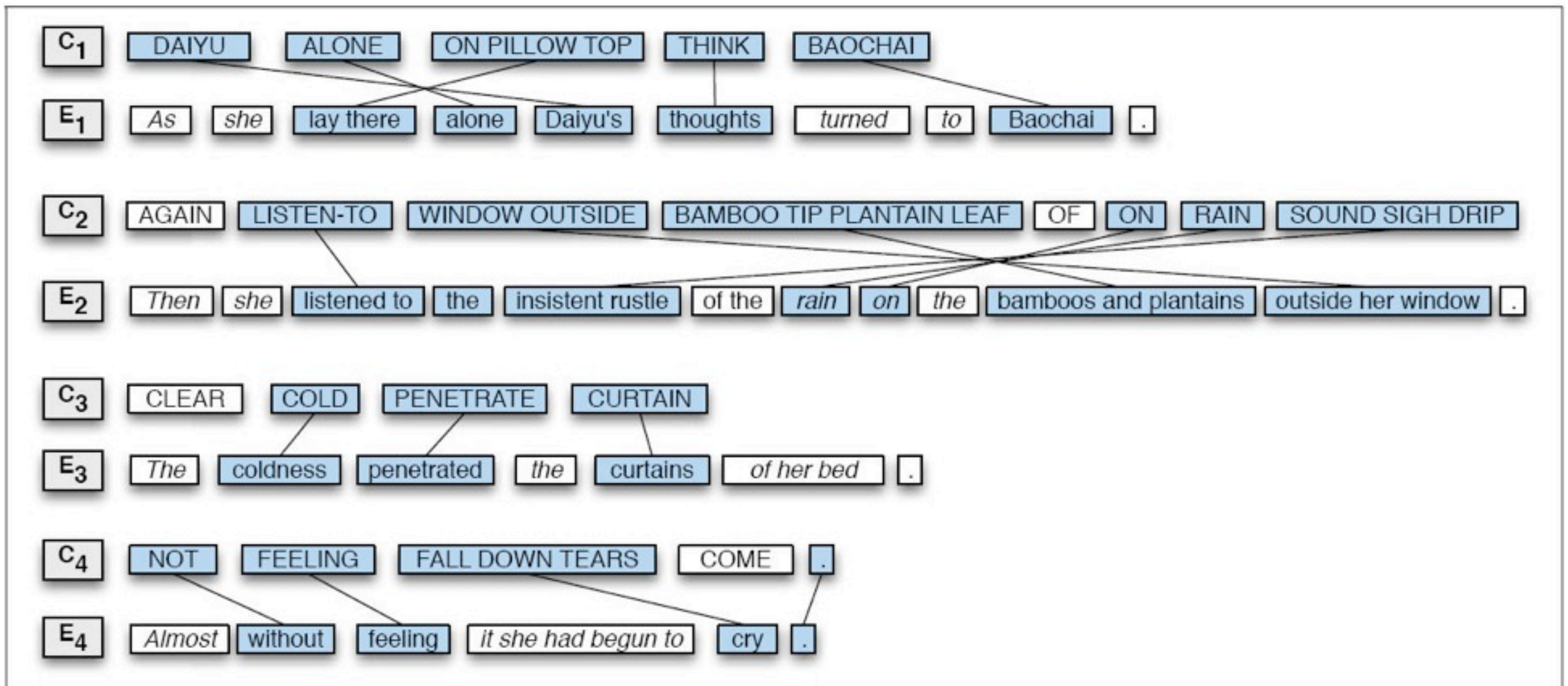
- Intro
- Classic MT
- Statistical MT

- Training
- Evaluation

MT is amazing

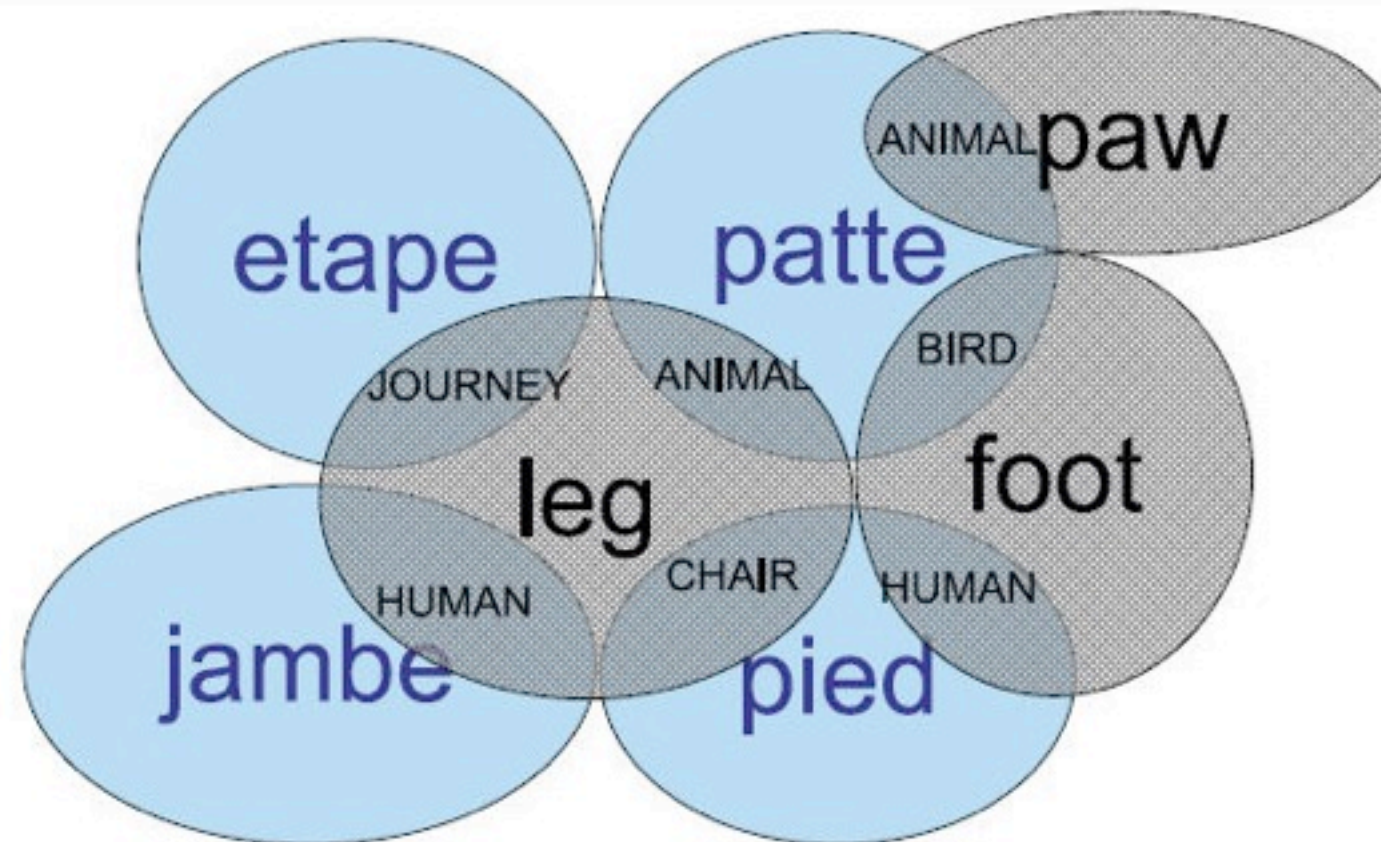
MT is hard

- Word order, word meanings



MT is hard

- Word meaning:
many-to-many and context dependent



- *Translation* itself is hard: metaphors, cultural references, etc.

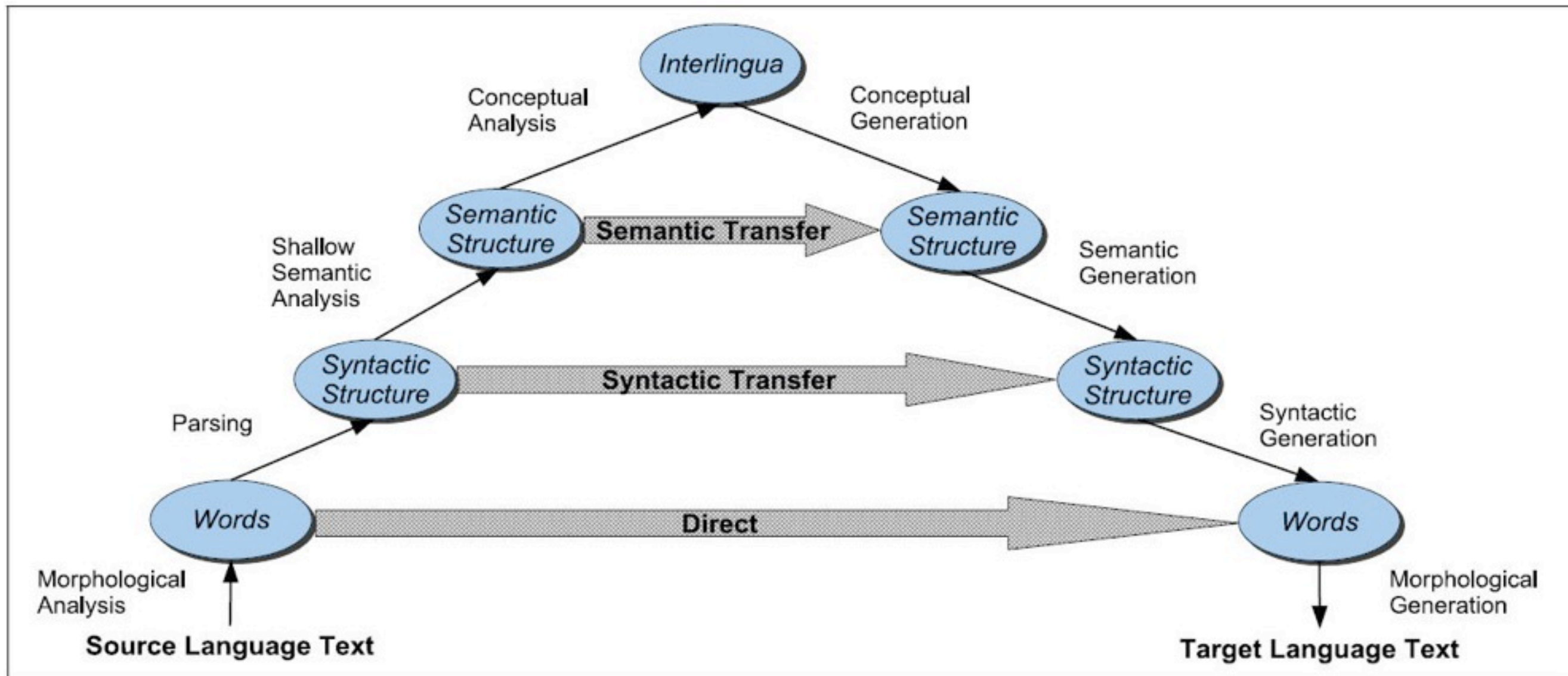
MT goals

- Motivation: Human translation is expensive
- High precision translation
- Rough translation
- Assistance for human translators
 - Comparison: bilingual dictionary

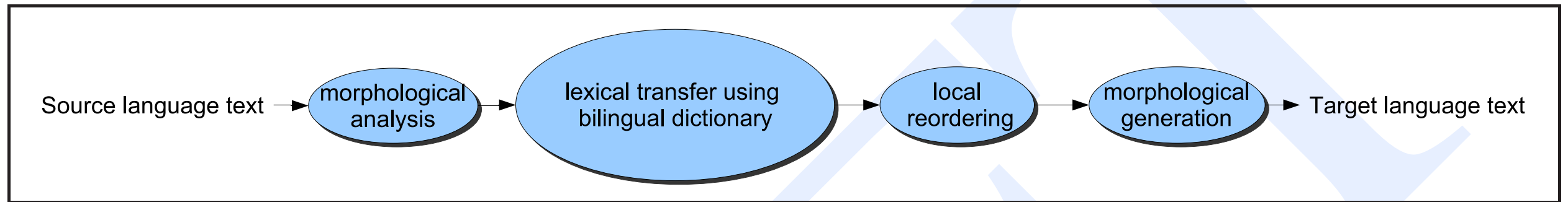
MT: major types

- Rule-based transfer
 - Manually program lexicons/rules
 - SYSTRAN (AltaVista Babelfish)
- Statistical MT:
 - Learn translation rules from data, search for high-scoring translation outputs
 - Phrase or syntactic transformations
 - Key research in the early 90s
 - Google Translate (mid 00s)
 - Moses, cdec (open-source)
- [Active current work: Semantic MT? Neural MT?]

Vauquois Triangle



Direct (word-based) transfer



Input:

After 1: Morphology

After 2: Lexical Transfer

After 3: Local reordering

After 4: Morphology

Mary didn't slap the green witch

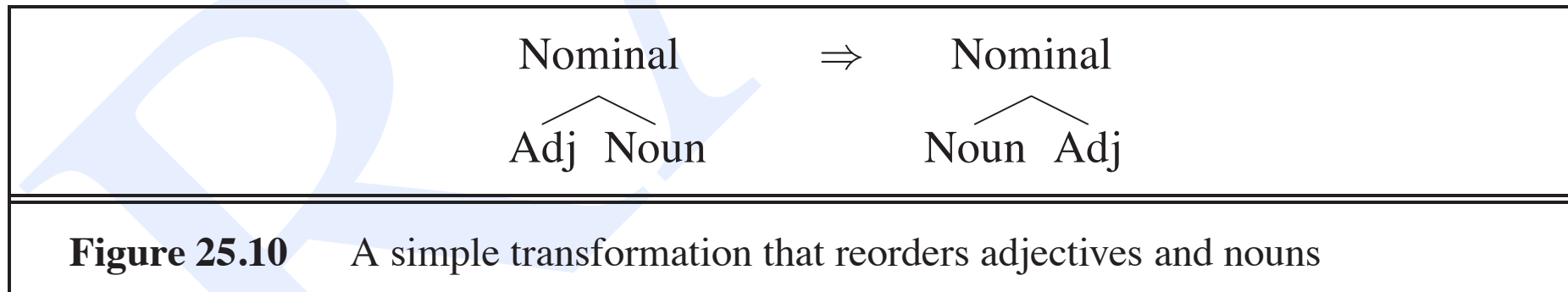
Mary DO-PAST not slap the green witch

Maria PAST no dar una bofetada a la verde bruja

Maria no dar PAST una bofetada a la bruja verde

Maria no dió una bofetada a la bruja verde

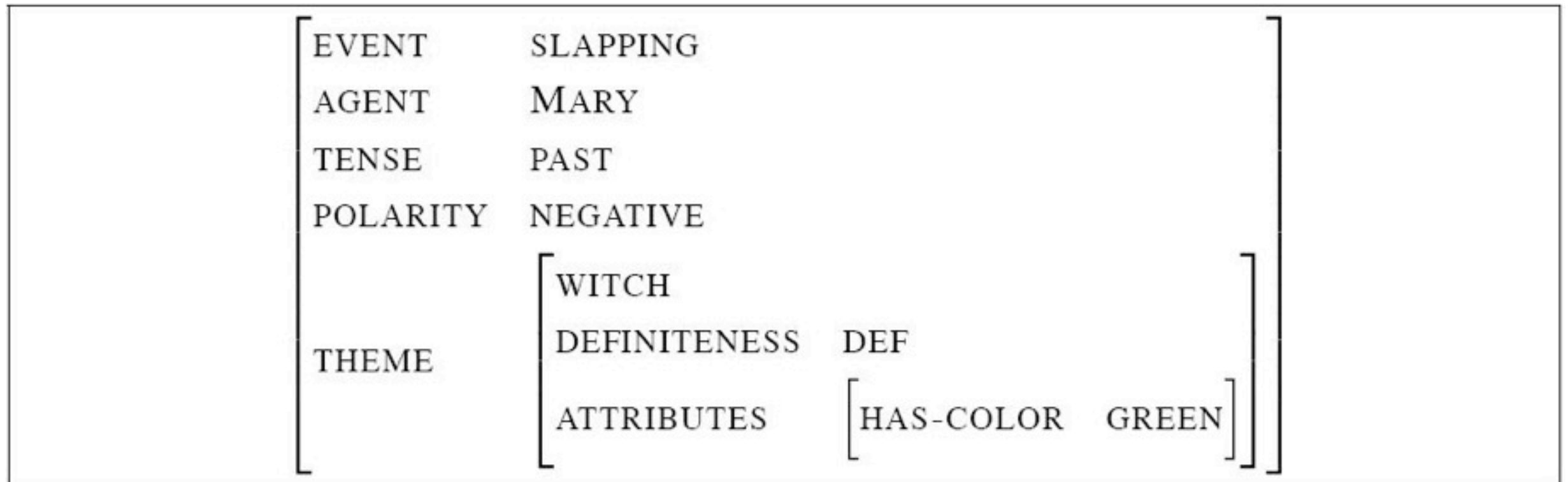
Syntactic transfer



English to Spanish:			
1.	$\text{NP} \rightarrow \text{Adjective}_1 \text{Noun}_2$	\Rightarrow	$\text{NP} \rightarrow \text{Noun}_2 \text{Adjective}_1$
Chinese to English:			
2.	$\text{VP} \rightarrow \text{PP}[+\text{Goal}] \text{V}$	\Rightarrow	$\text{VP} \rightarrow \text{V PP}[+\text{Goal}]$
English to Japanese:			
3.	$\text{VP} \rightarrow \text{V NP}$	\Rightarrow	$\text{VP} \rightarrow \text{NP V}$
4.	$\text{PP} \rightarrow \text{P NP}$	\Rightarrow	$\text{PP} \rightarrow \text{NP P}$
5.	$\text{NP} \rightarrow \text{NP}_1 \text{Rel. Clause}_2$	\Rightarrow	$\text{NP} \rightarrow \text{Rel. Clause}_2 \text{NP}_1$

Interlingua

“Mary did not slap the green witch”



- More like classic logic-based AI
- Works in narrow domains
- Broad domain currently fails
 - Coverage: Knowledge representation for all possible semantics?
 - Can you parse to it?
 - Can you generate from it?

Rules are hard

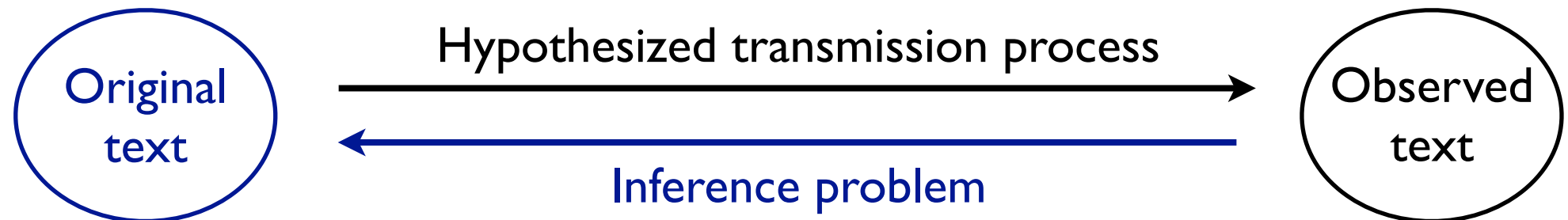
- Coverage
- Complexity (context dependence)
- Maintenance

```
function DIRECT_TRANSLATE_MUCH/MANY(word) returns Russian translation  
  
if preceding word is how return skol'ko  
else if preceding word is as return stol'ko zhe  
else if word is much  
    if preceding word is very return nil  
    else if following word is a noun return mnogo  
else /* word is many */  
    if preceding word is a preposition and following word is a noun return mnogii  
    else return mnogo
```

Statistical MT

- MT as ML: Translation is something people do naturally. Learn rules from data?
- Parallel data: (source, target) text pairs
 - E.g. 20 million words of European Parliament proceedings
<http://www.statmt.org/euoparl/>

Noisy channel model



Codebreaking

$$P(\text{plaintext} \mid \text{encrypted text}) \propto P(\text{encrypted text} \mid \text{plaintext}) P(\text{plaintext})$$

Speech recognition

$$P(\text{text} \mid \text{acoustic signal}) \propto P(\text{acoustic signal} \mid \text{text}) P(\text{text})$$

Optical character recognition

$$P(\text{text} \mid \text{image}) \propto P(\text{image} \mid \text{text}) P(\text{text})$$

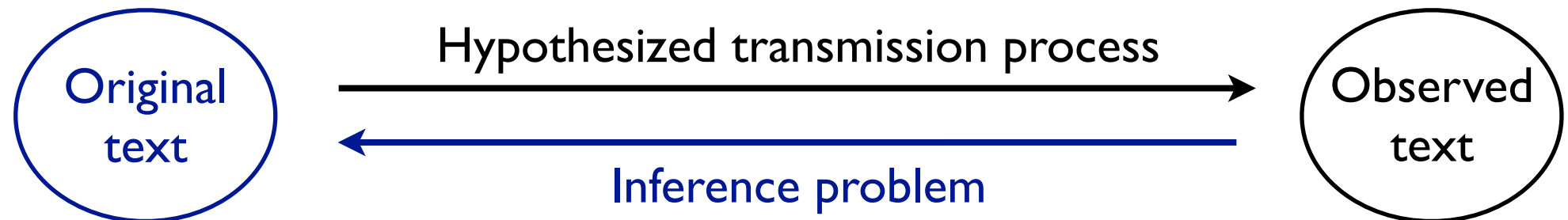
Machine translation

$$P(\text{target text} \mid \text{source text}) \propto P(\text{source text} \mid \text{target text}) P(\text{target text})$$

Spelling correction

$$P(\text{target text} \mid \text{source text}) \propto P(\text{source text} \mid \text{target text}) P(\text{target text})$$

Noisy channel model



Codebreaking

$$P(\text{plaintext} | \text{ciphertext}) \propto P(\text{ciphertext} | \text{plaintext}) P(\text{plaintext})$$

Speech

$$P(\text{text} | \text{audio})$$

Optical

$$P(\text{text} | \text{image})$$

Machine translation

$$P(\text{target text} | \text{source text}) \propto P(\text{source text} | \text{target text}) P(\text{target text})$$

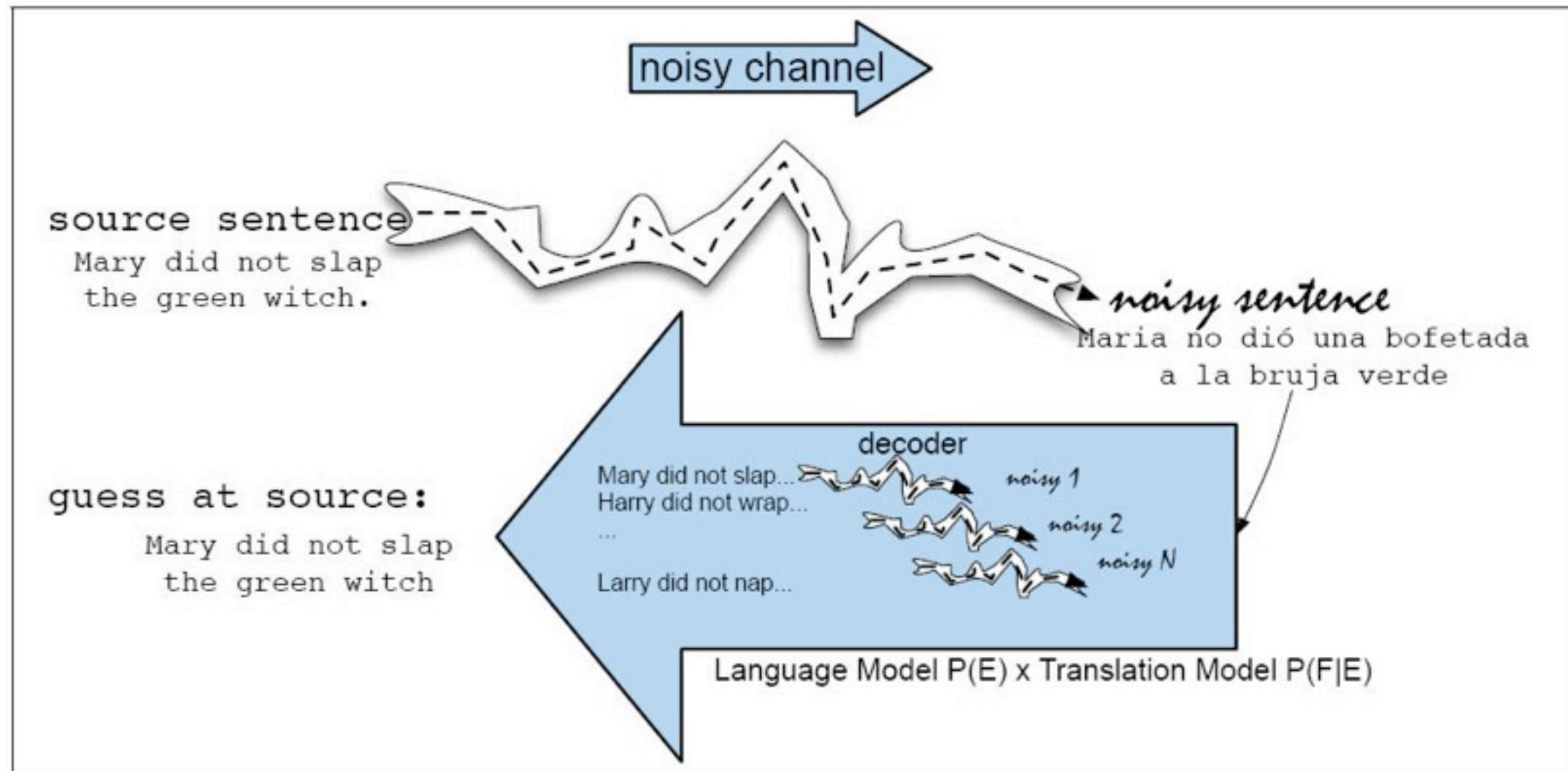
Spelling correction

$$P(\text{target text} | \text{source text}) \propto P(\text{source text} | \text{target text}) P(\text{target text})$$

One naturally wonders if the problem of translation could conceivably be treated as a problem in cryptography. 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 (1955)

Statistical MT



Statistical MT

- Pioneered at IBM, early 1990s
(Forerunner of 90s-era statistical revolution in NLP)

The COLING Paper Review

The validity of statistical (information theoretic) approach to MT has indeed been recognized, as the authors mention, by Weaver as early as 1949. And was universally recognized as mistaken by 1950. (cf. Hutchins, MT: Past, Present, Future, Ellis Horwood, 1986, pp. 30ff. and references therein) The crude force of computers is not science. The paper is simply beyond the scope of COLING.

Historical notes: <http://cs.jhu.edu/~post/bitext/>

Statistical MT

- Pioneered at IBM, early 1990s
(Forerunner of 90s-era statistical revolution in NLP)
- Noisy channel model borrowed from
speech recognition processing


"Every time I fire a linguist,
the performance of the speech recognizer goes up"
[Fred Jelinek]

Problem formulation

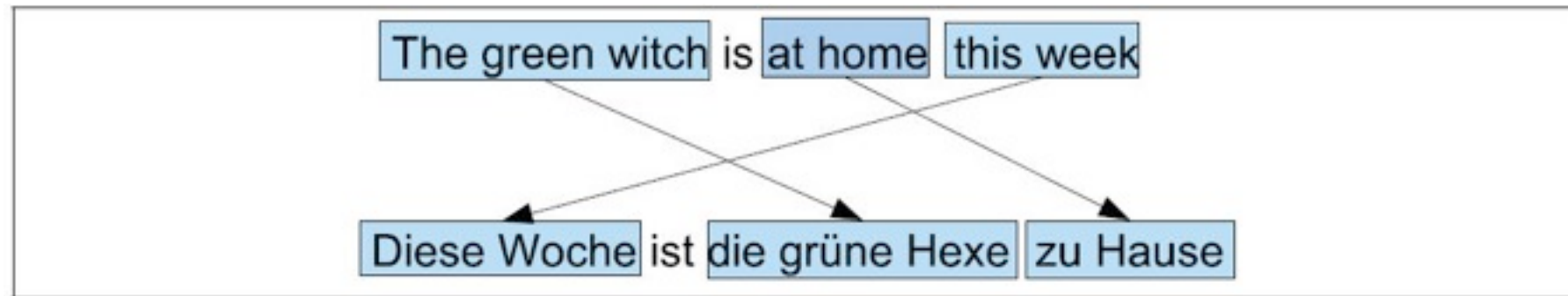
best-translation $\hat{T} = \operatorname{argmax}_T \text{faithfulness}(T,S) \text{ fluency}(T)$

$$\hat{E} = \operatorname{argmax}_{E \in \text{English}} \underbrace{P(F|E)}_{\text{translation model}} \underbrace{P(E)}_{\text{language model}}$$

↑
decoder
(search algo)



Phrase-based model



- Learning $P(F | E)$ phrase translation tables:
Assume aligned corpus. Then count
- Today: lexical translation model (IBM Model 1)

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


Maximum Likelihood Estimation: count ratios

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

Could learn ***if*** we had translation frequencies.

Lexical Translation

- Goal: a model $p(\mathbf{e} \mid \mathbf{f}, m)$
- where \mathbf{e} and \mathbf{f} are complete English and Foreign sentences

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


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

$$\mathbf{e} = \langle e_1, e_2, \dots, e_m \rangle \quad \mathbf{f} = \langle f_1, f_2, \dots, f_n \rangle$$
$$\mathbf{a} = \langle a_1, a_2, \dots, a_m \rangle \quad \text{each } a_i \in \{0, 1, \dots, n\}$$

Modeling assumptions

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

[Alignment] × [Translation | Alignment]

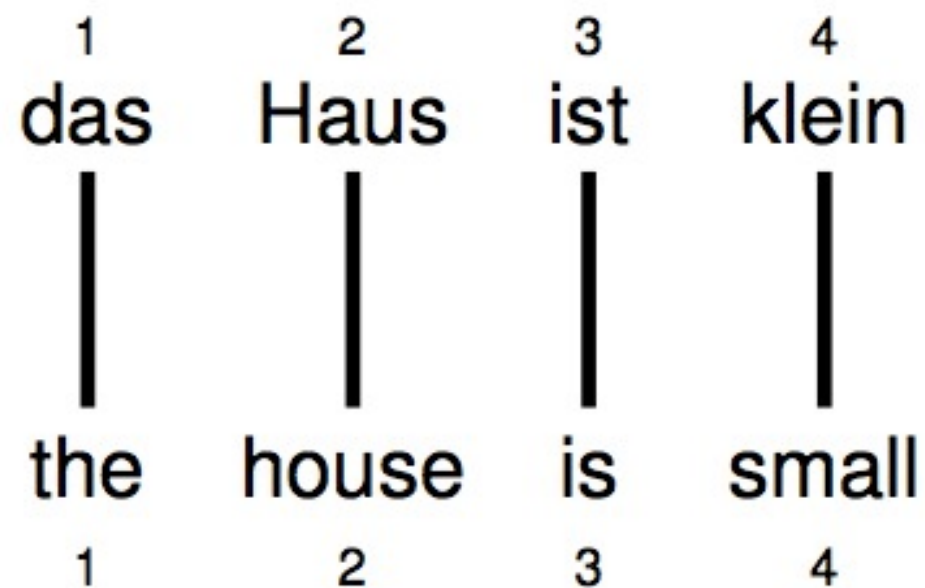
Alignment

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

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

Alignment

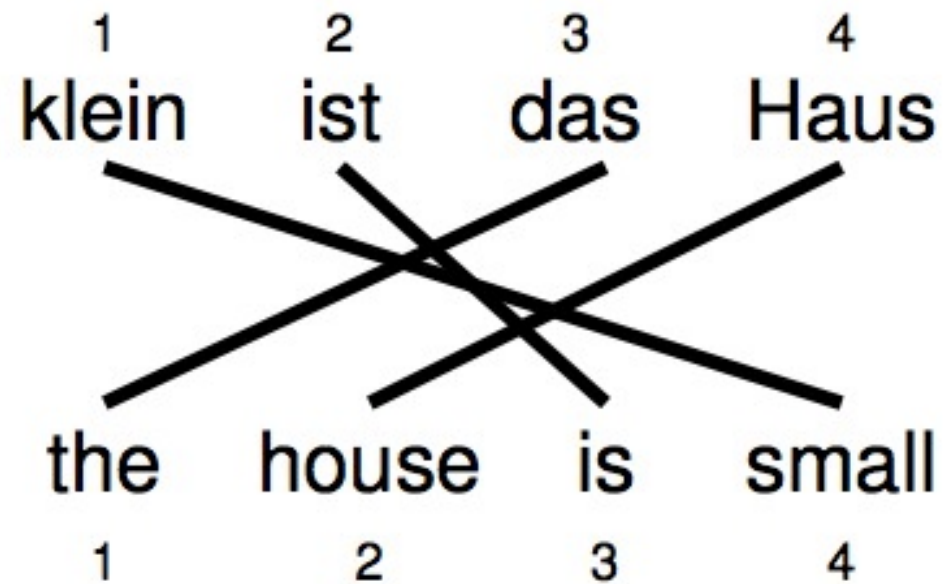
- Alignments can be visualized in by drawing links between two sentences, and they are represented as vectors of positions:



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

Reordering

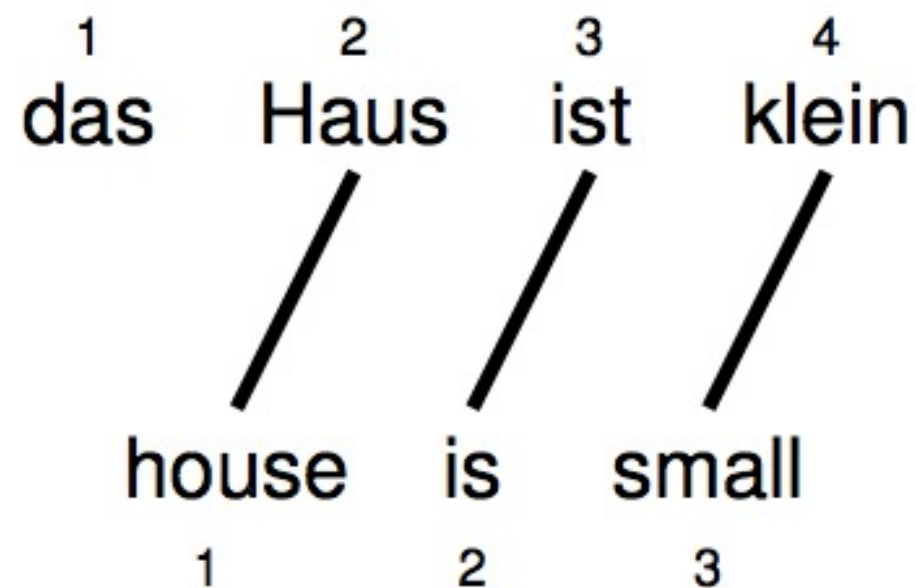
- Words may be reordered during translation.



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

Word Dropping

- A source word may not be translated at all



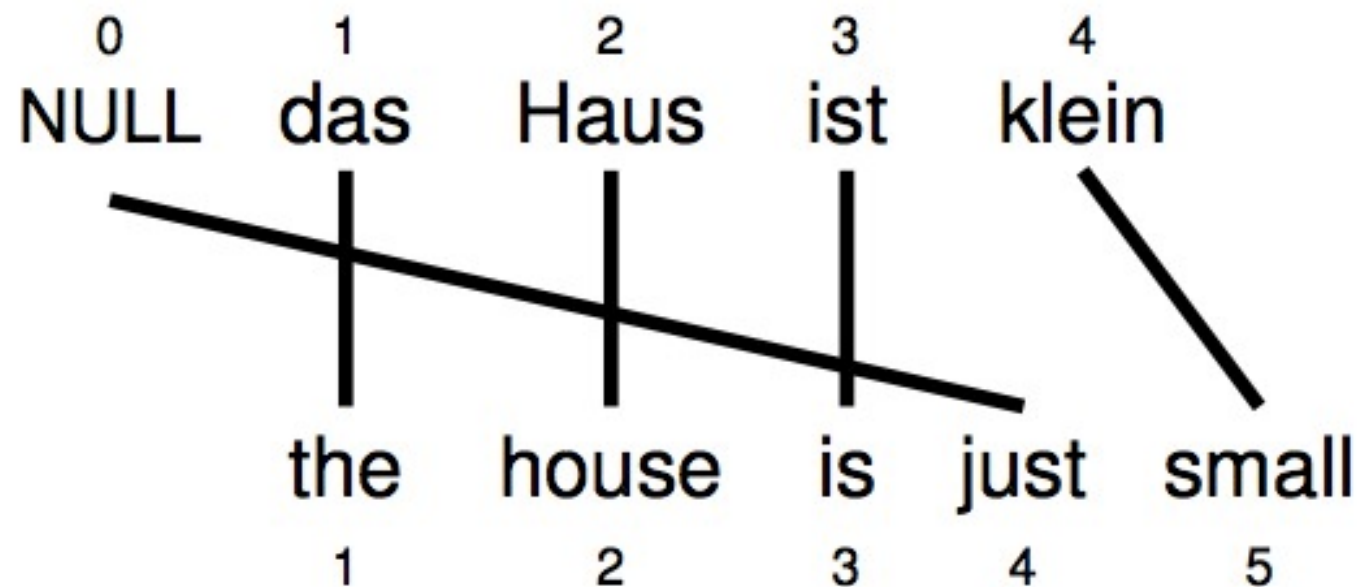
$$\mathbf{a} = (2, 3, 4)^T$$

Word Insertion

- Words may be inserted during translation

English **just** does not have an equivalent

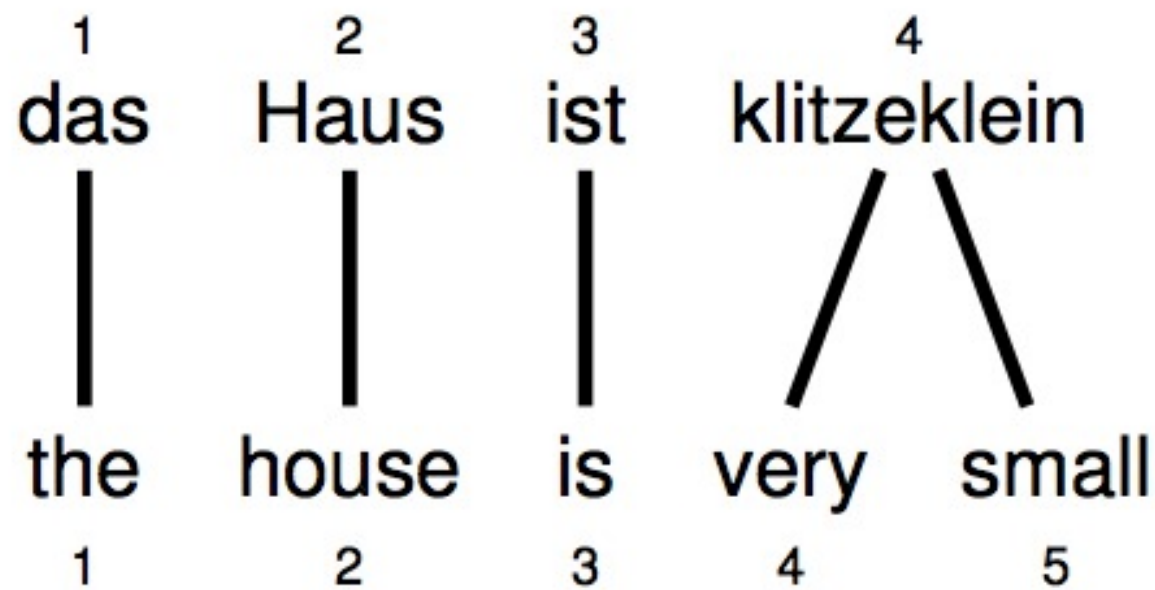
But it must be explained - we typically assume every source sentence contains a NULL token



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

One-to-many Translation

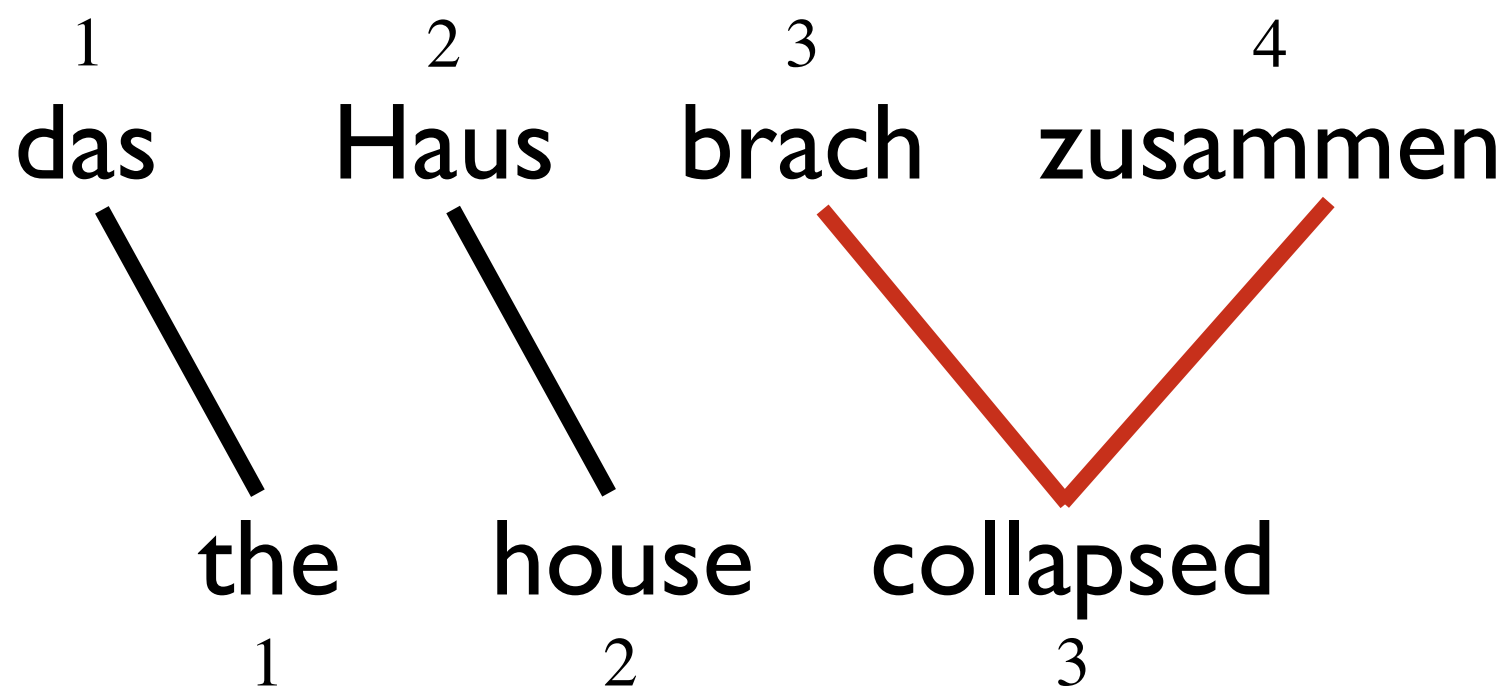
- A source word may translate into **more than one** target word



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

Many-to-one Translation

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



$a = ???$

[IBM Model 1 can't do this]

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

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

$$e_i \sim \text{Categorical}(\boldsymbol{\theta}_{f_{a_i}})$$

IBM Model I

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

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

$e_i \sim \text{Categorical}(\boldsymbol{\theta}_{f_{a_i}})$

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

IBM Model I

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

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

$$e_i \sim \text{Categorical}(\boldsymbol{\theta}_{f_{a_i}})$$

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

IBM Model I

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

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

$e_i \sim \text{Categorical}(\boldsymbol{\theta}_{f_{a_i}})$

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

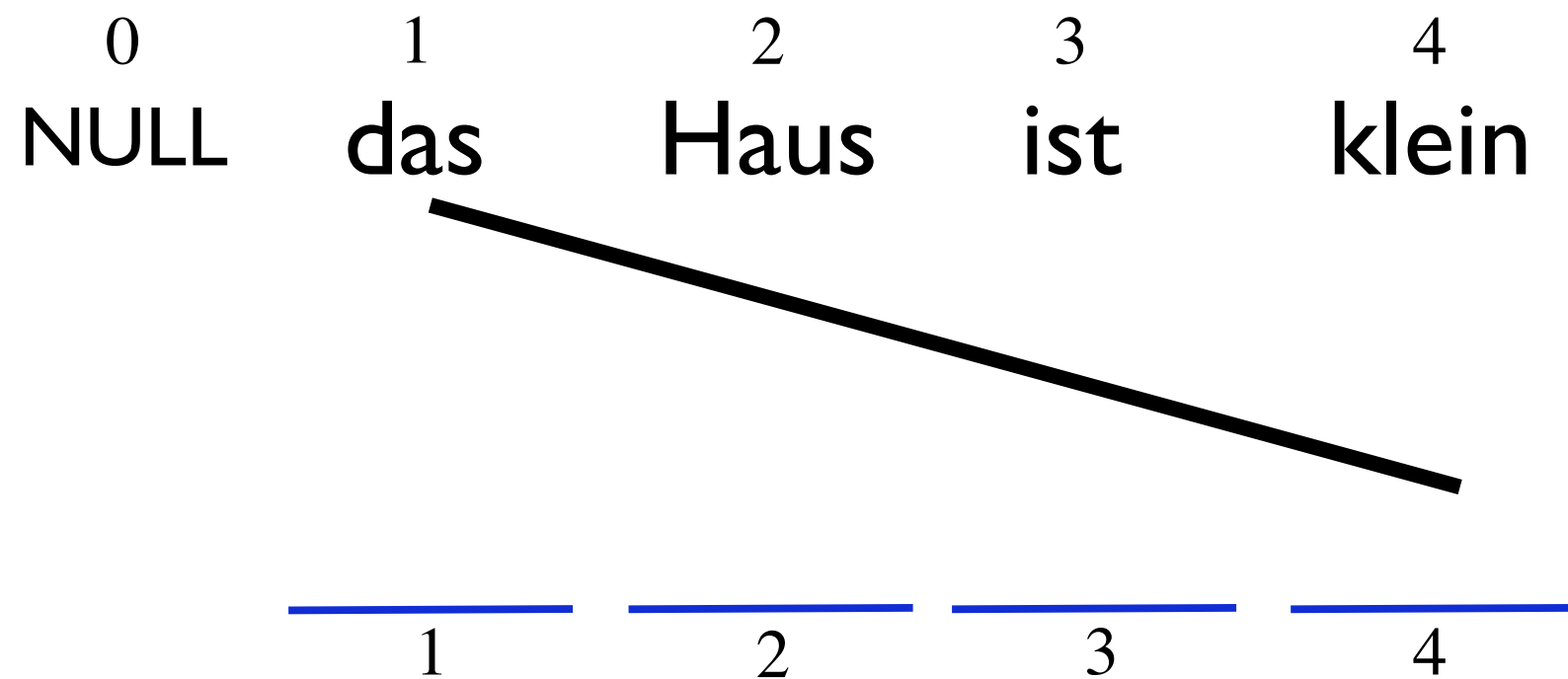
Example

0 1 2 3 4
NULL das Haus ist klein

— — — —
1 2 3 4

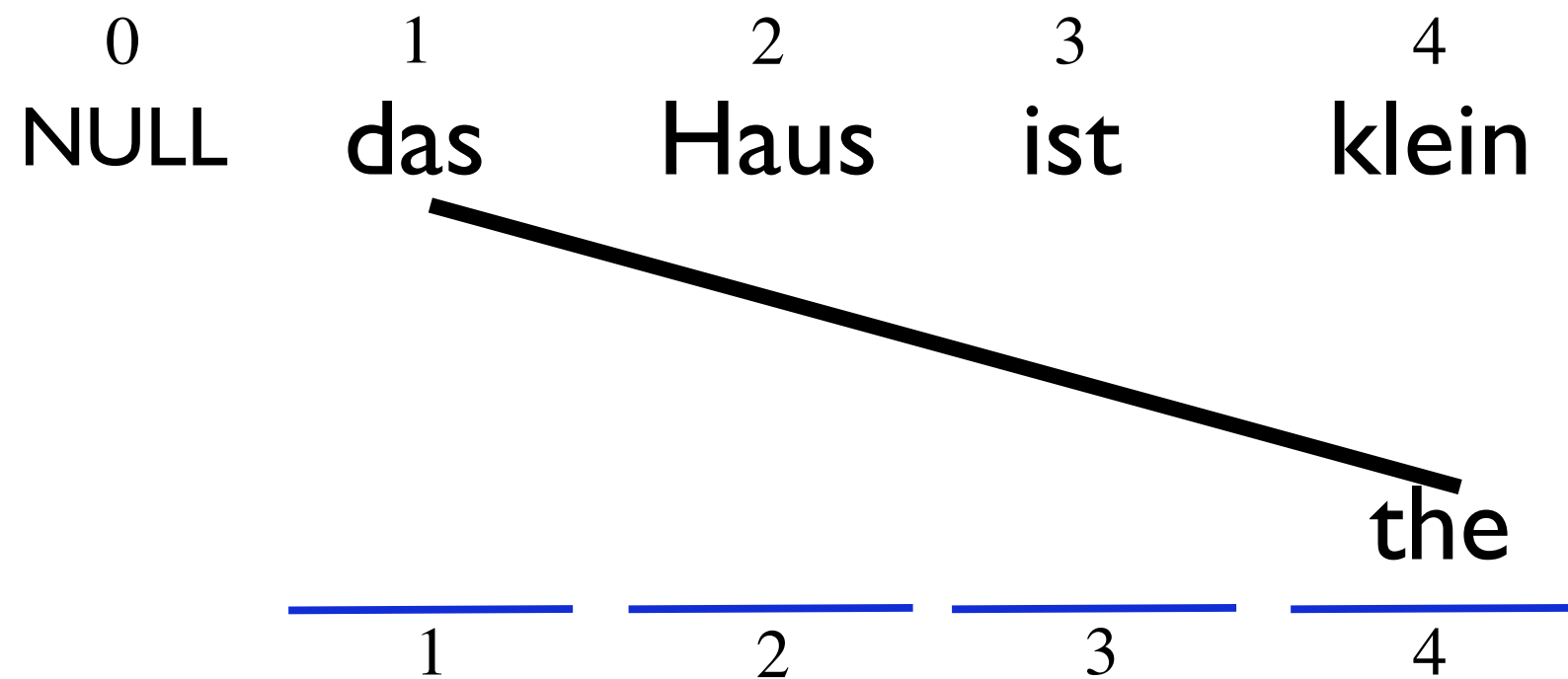
$p(\mathbf{e} \mid \mathbf{f})$: Assume a foreign sentence and target length.
 $= p(\mathbf{a} \mid \mathbf{f}) p(\mathbf{e} \mid \mathbf{a}, \mathbf{f})$

Example



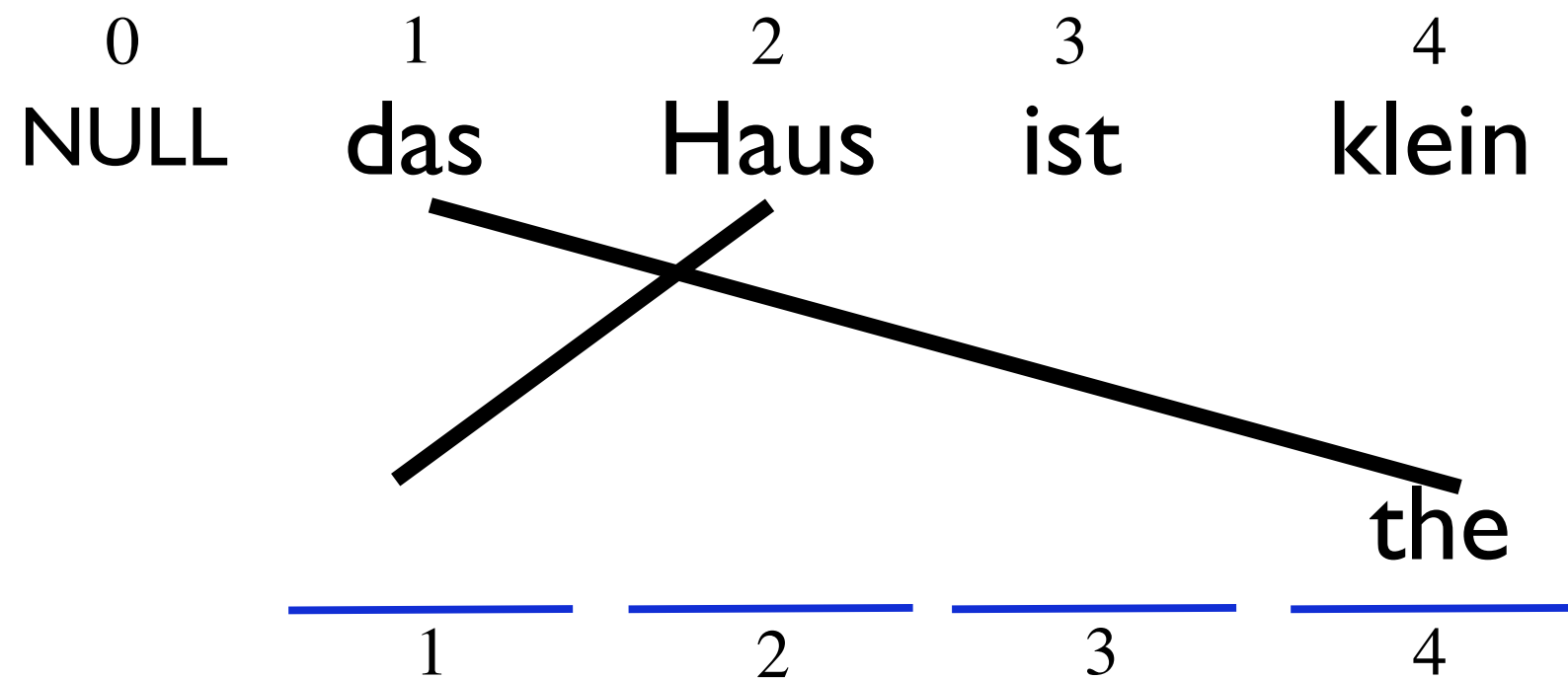
$p(\mathbf{e} \mid \mathbf{f})$: Assume a foreign sentence and target length.
 $= p(\mathbf{a} \mid \mathbf{f}) p(\mathbf{e} \mid \mathbf{a}, \mathbf{f})$

Example



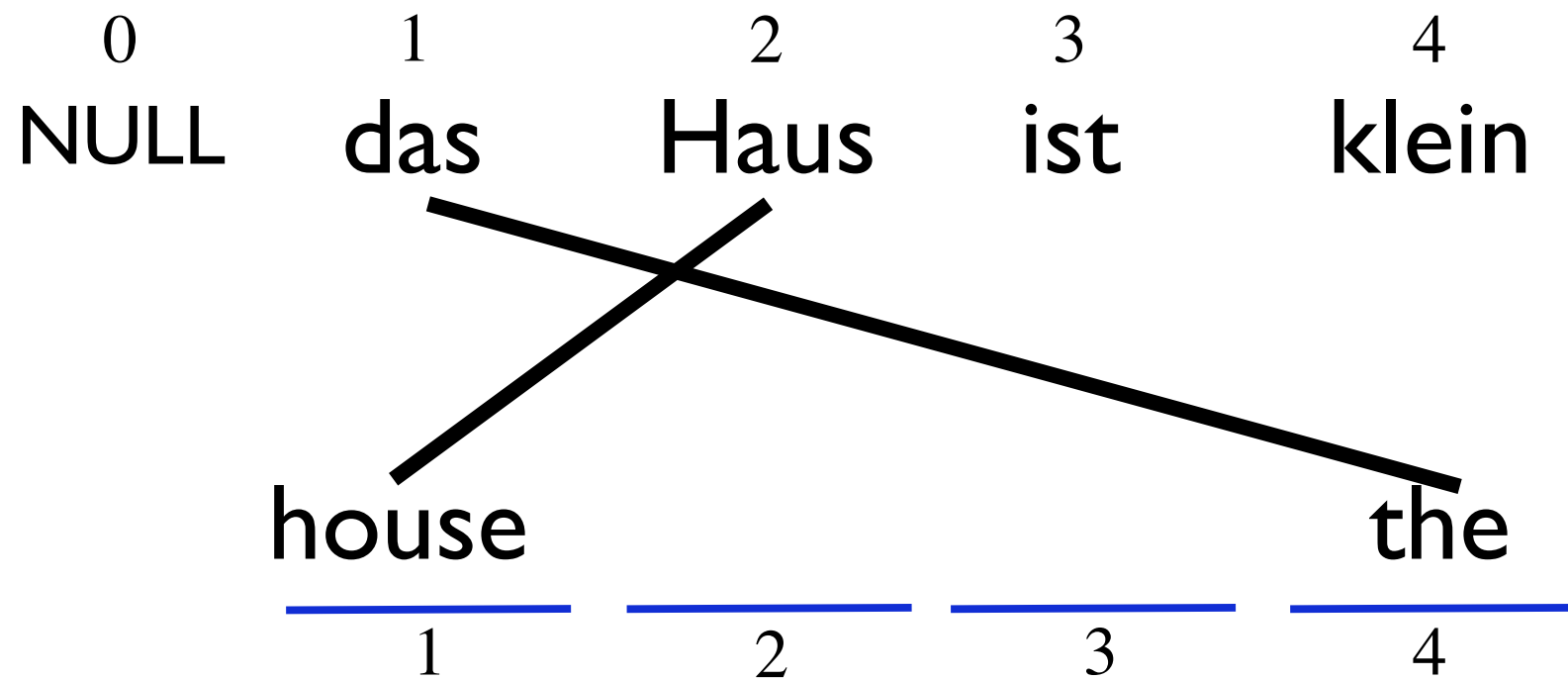
$p(\mathbf{e} \mid \mathbf{f})$: Assume a foreign sentence and target length.
 $= p(\mathbf{a} \mid \mathbf{f}) p(\mathbf{e} \mid \mathbf{a}, \mathbf{f})$

Example



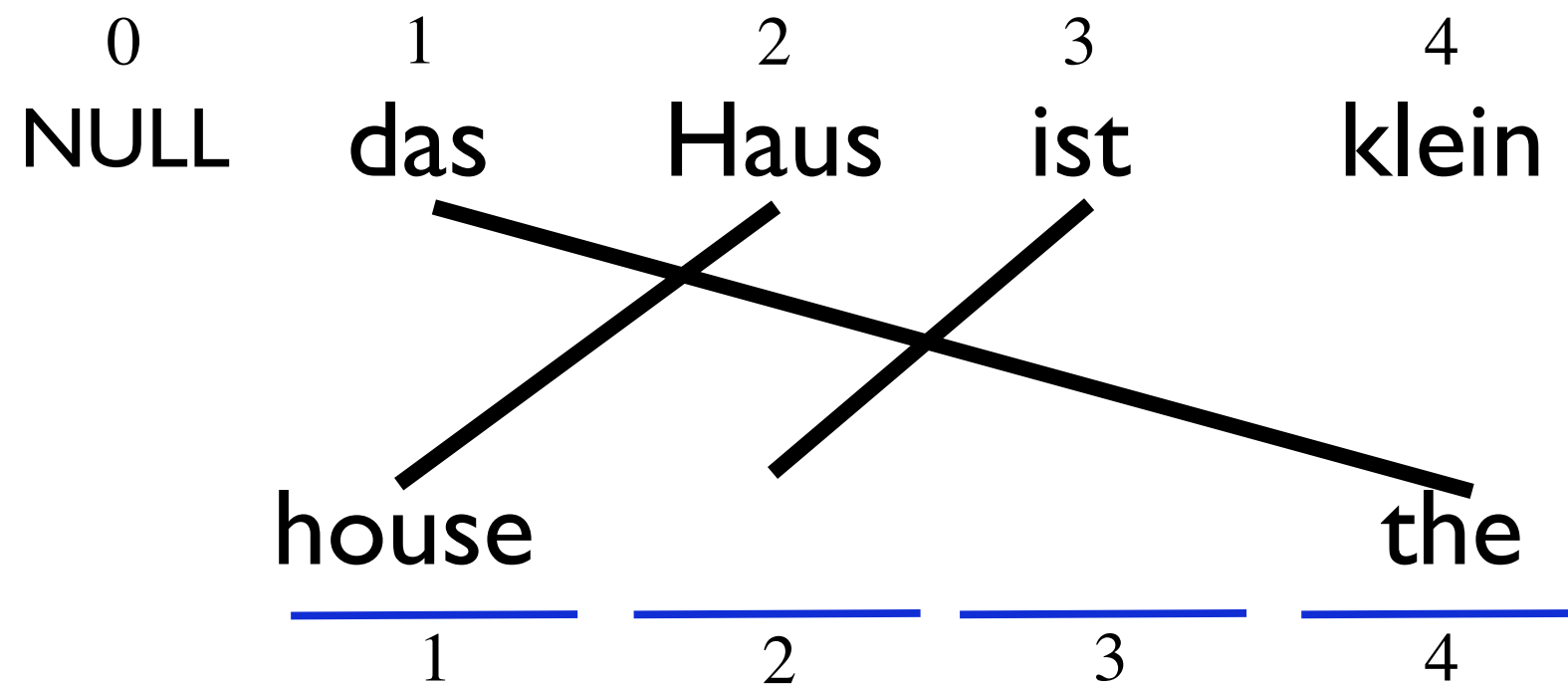
$p(\mathbf{e} \mid \mathbf{f})$: Assume a foreign sentence and target length.
 $= p(\mathbf{a} \mid \mathbf{f}) p(\mathbf{e} \mid \mathbf{a}, \mathbf{f})$

Example



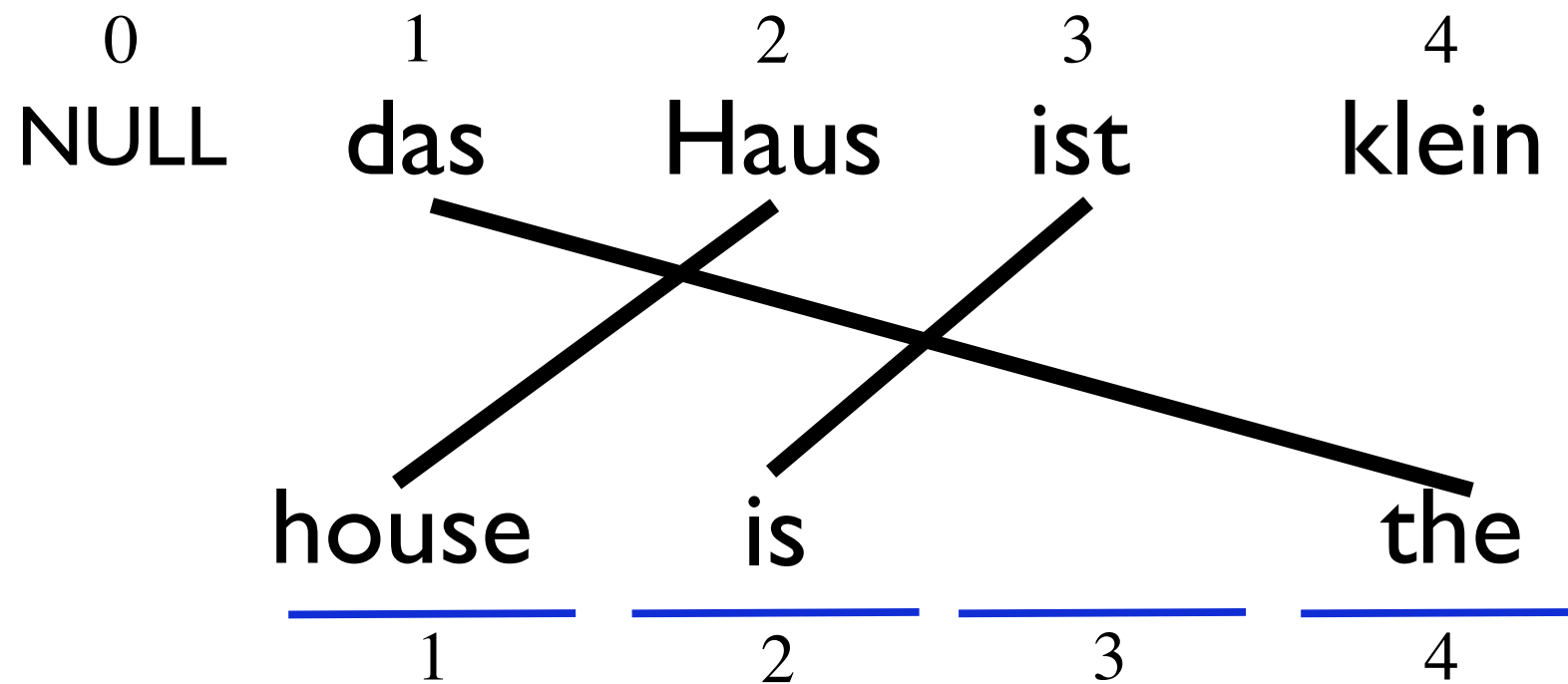
$p(\mathbf{e} \mid \mathbf{f})$: Assume a foreign sentence and target length.
 $= p(\mathbf{a} \mid \mathbf{f}) p(\mathbf{e} \mid \mathbf{a}, \mathbf{f})$

Example



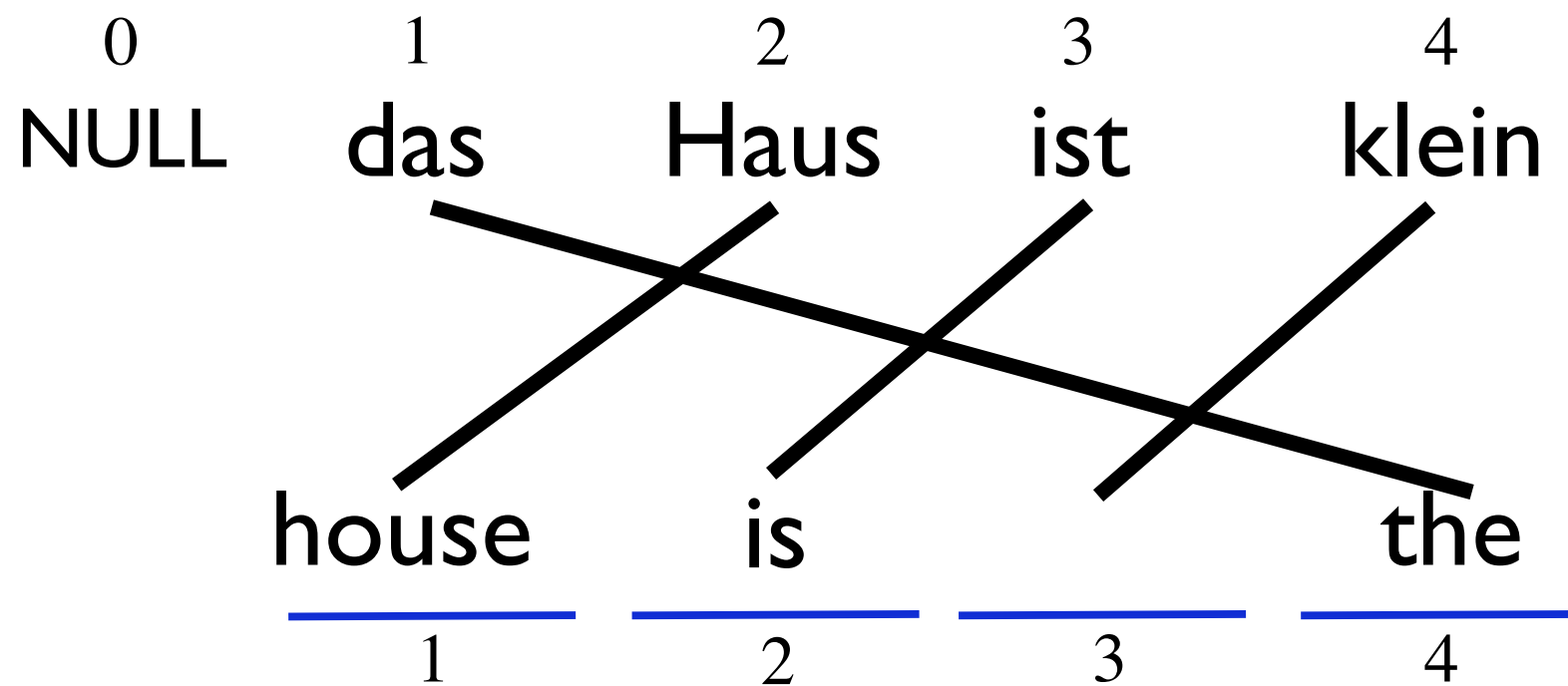
$p(\mathbf{e} \mid \mathbf{f})$: Assume a foreign sentence and target length.
 $= p(\mathbf{a} \mid \mathbf{f}) p(\mathbf{e} \mid \mathbf{a}, \mathbf{f})$

Example



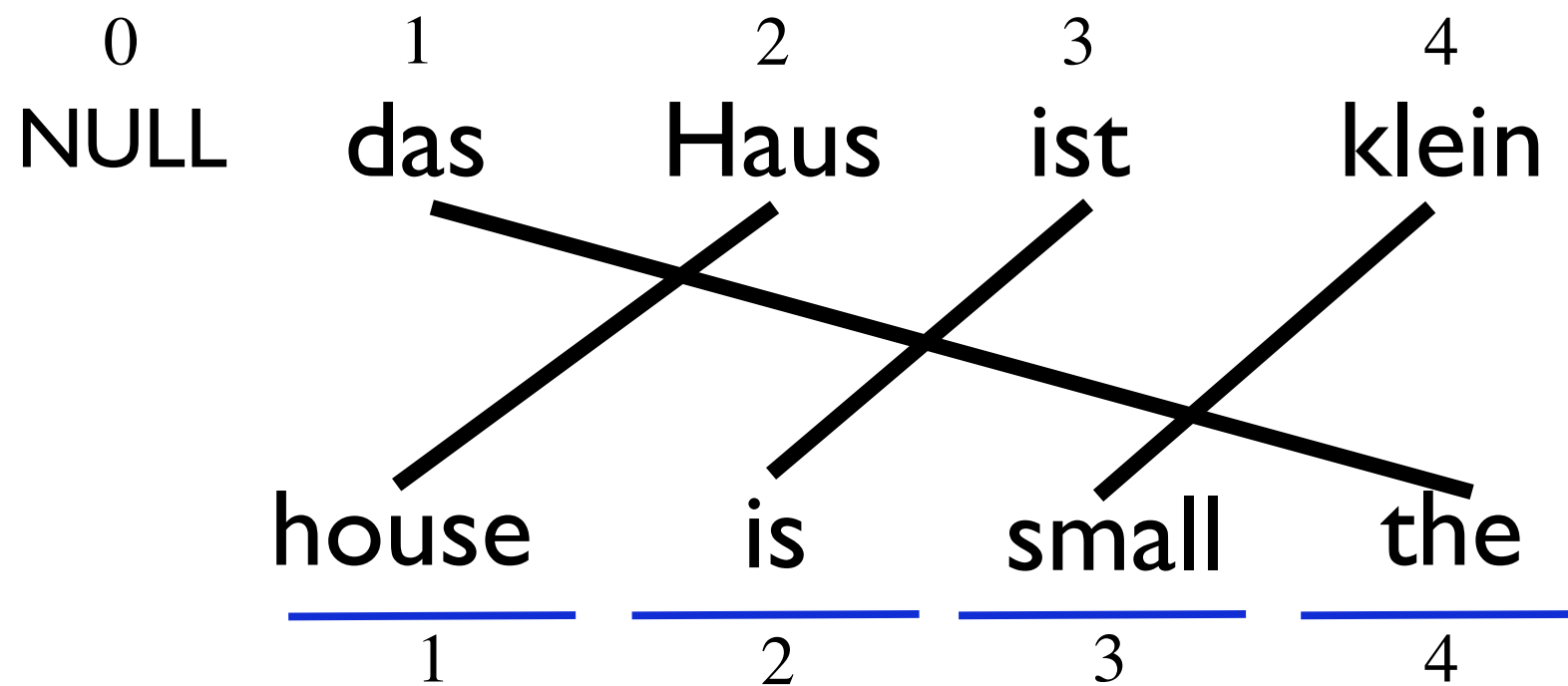
$p(\mathbf{e} \mid \mathbf{f})$: Assume a foreign sentence and target length.
 $= p(\mathbf{a} \mid \mathbf{f}) p(\mathbf{e} \mid \mathbf{a}, \mathbf{f})$

Example



$p(\mathbf{e} \mid \mathbf{f})$: Assume a foreign sentence and target length.
 $= p(\mathbf{a} \mid \mathbf{f}) p(\mathbf{e} \mid \mathbf{a}, \mathbf{f})$

Example



$p(\mathbf{e} \mid \mathbf{f})$: Assume a foreign sentence and target length.
 $= p(\mathbf{a} \mid \mathbf{f}) p(\mathbf{e} \mid \mathbf{a}, \mathbf{f})$

IBM Model I: Inference and learning

- Alignment inference:
Given lexical translation probabilities,
infer posterior or Viterbi alignment

$$\arg \max_{\mathbf{a}} p(\mathbf{a} \mid \mathbf{e}, \mathbf{f}, \theta)$$

- Translation: incorporate into noisy channel
(this model isn't good at this)

$$\arg \max_e p(\mathbf{e} \mid \mathbf{f}, \theta) p(\mathbf{e})$$

- How do we learn translation parameters?
EM Algorithm (Thursday)

$$\arg \max_{\theta} p(\mathbf{e} \mid \mathbf{f}, \theta)$$

- Chicken and egg problem:
If we knew alignments, translation
parameters would be trivial (just counting)

