

Lexical Translation Models I

Machine Translation Lecture 4

Instructor: Chris Callison-Burch
TAs: Mitchell Stern, Justin Chiu

Website: mt-class.org/penn



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

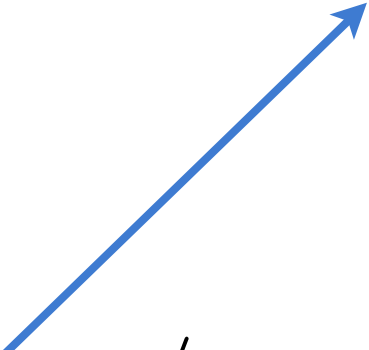
$$\hat{p}_{\text{MLE}}(e \mid \text{Haus}) = \begin{cases} 0.696 & \text{if } e = \text{house} \\ 0.279 & \text{if } e = \text{home} \\ 0.014 & \text{if } e = \text{shell} \\ 0.011 & \text{if } e = \text{household} \\ 0 & \text{otherwise} \end{cases}$$

Lexical Translation

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

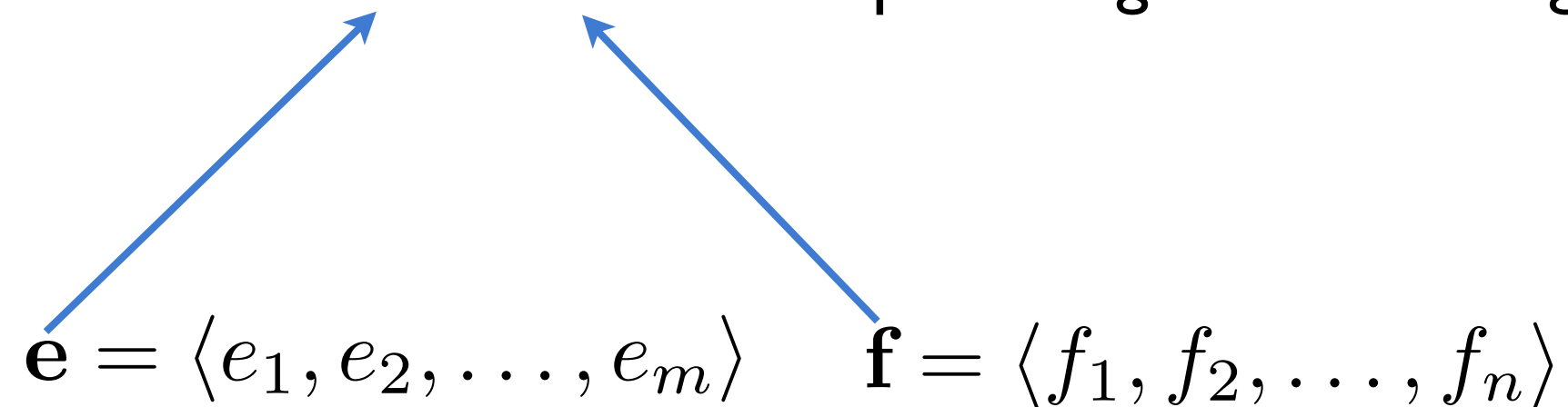
Lexical Translation

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

$$\mathbf{e} = \langle e_1, e_2, \dots, e_m \rangle$$


Lexical Translation

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



The diagram consists of two blue arrows pointing upwards from the definitions of \mathbf{e} and \mathbf{f} to the variables \mathbf{e} and \mathbf{f} in the list above. The first arrow starts at the \mathbf{e} in the equation $\mathbf{e} = \langle e_1, e_2, \dots, e_m \rangle$ and points to the \mathbf{e} in the text "where \mathbf{e} and \mathbf{f} are complete English and Foreign sentences". The second arrow starts at the \mathbf{f} in the equation $\mathbf{f} = \langle f_1, f_2, \dots, f_n \rangle$ and points to the \mathbf{f} in the same text.

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

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 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:

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

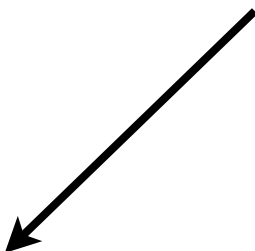
Alignment \times Translation | Alignment

Lexical Translation

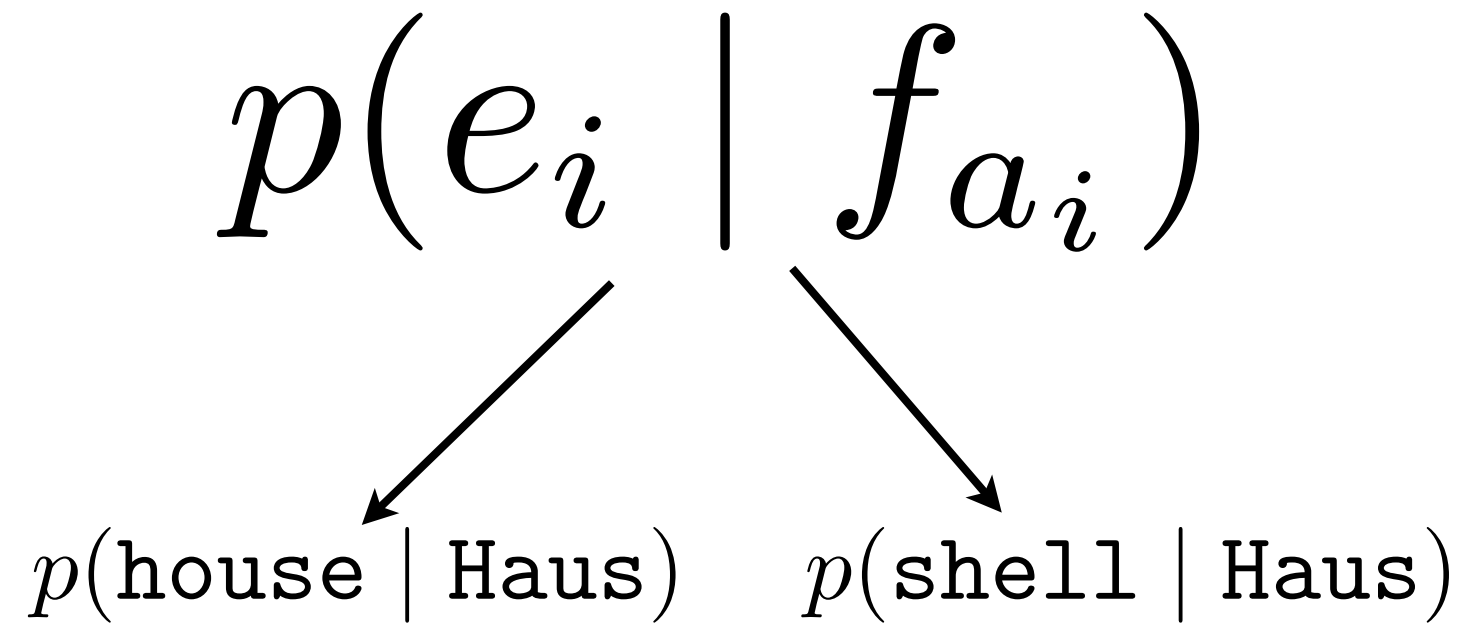
$$p(e_i \mid f_{a_i})$$

Lexical Translation

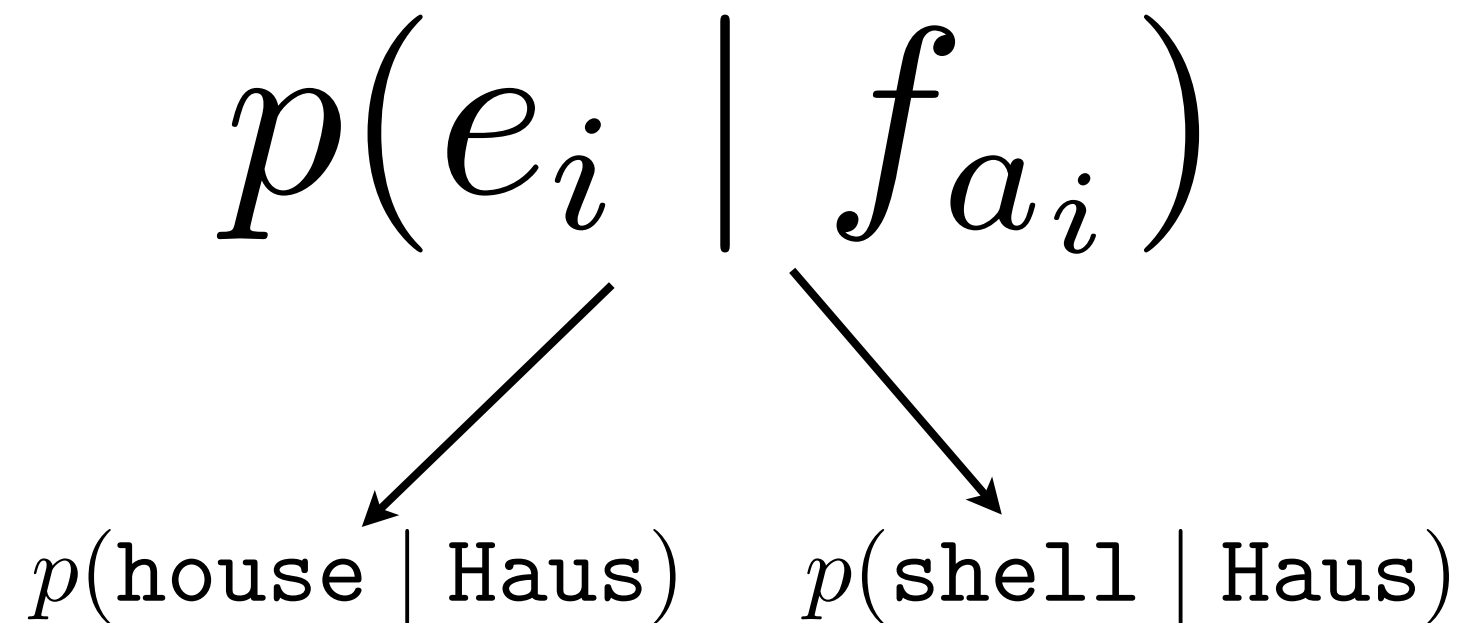
$$p(e_i \mid f_{a_i})$$


$$p(\text{house} \mid \text{Haus})$$

Lexical Translation



Lexical Translation



Remember bigram models...

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

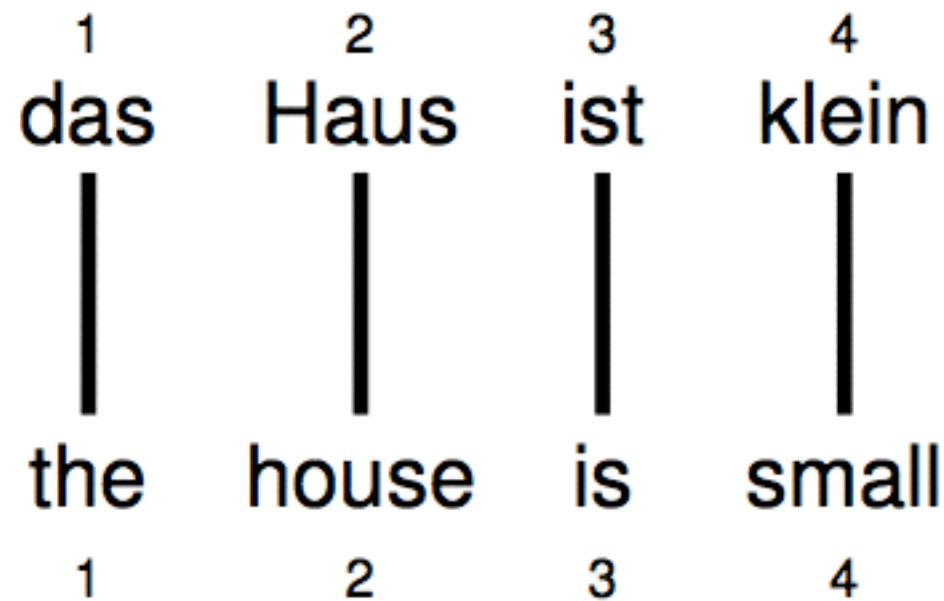
Alignment

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

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

Alignment

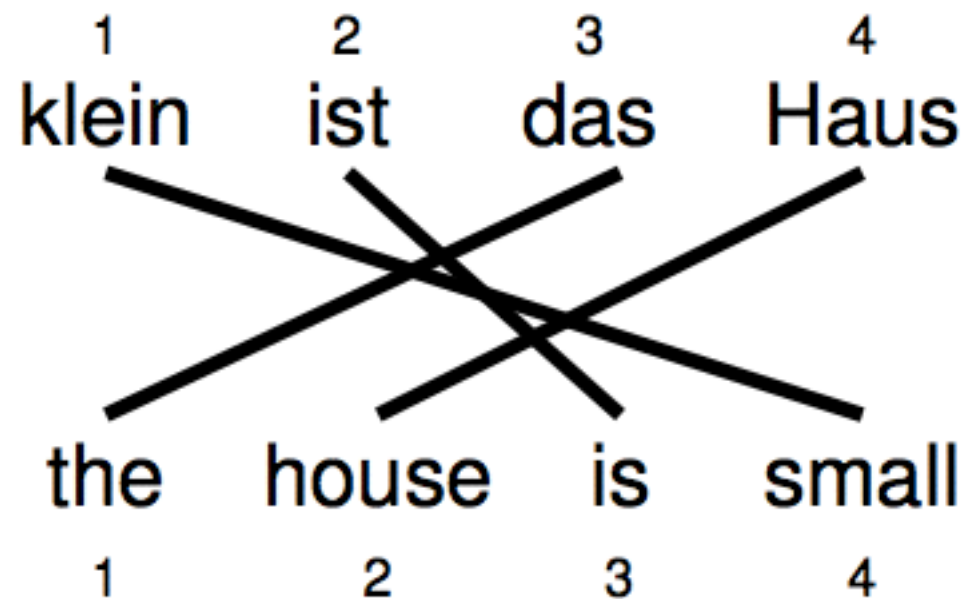
- Alignments can be visualized in by drawing links between two sentences, and they are represented as vectors of positions:



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

Reordering

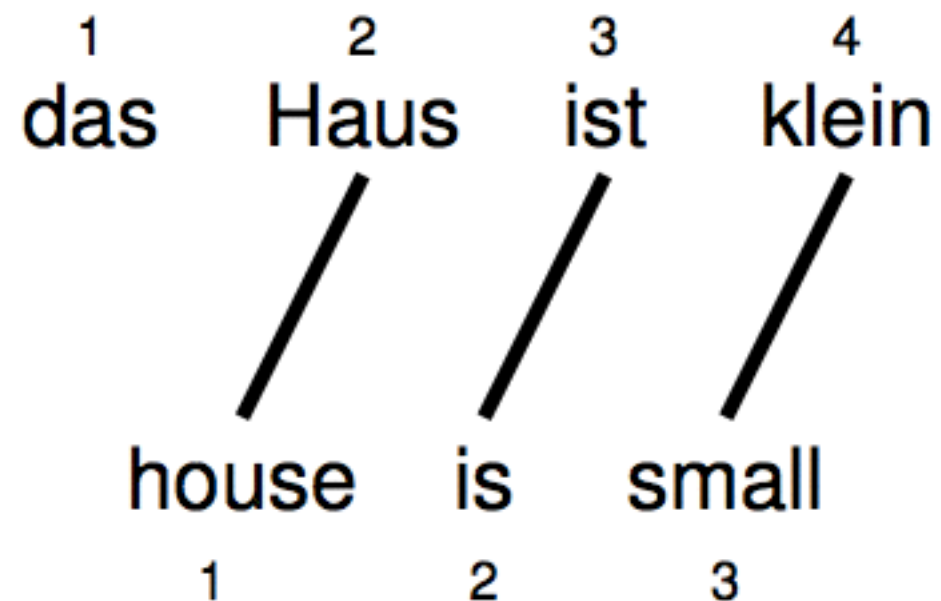
- Words may be reordered during translation.



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

Word Dropping

- A source word may not be translated at all

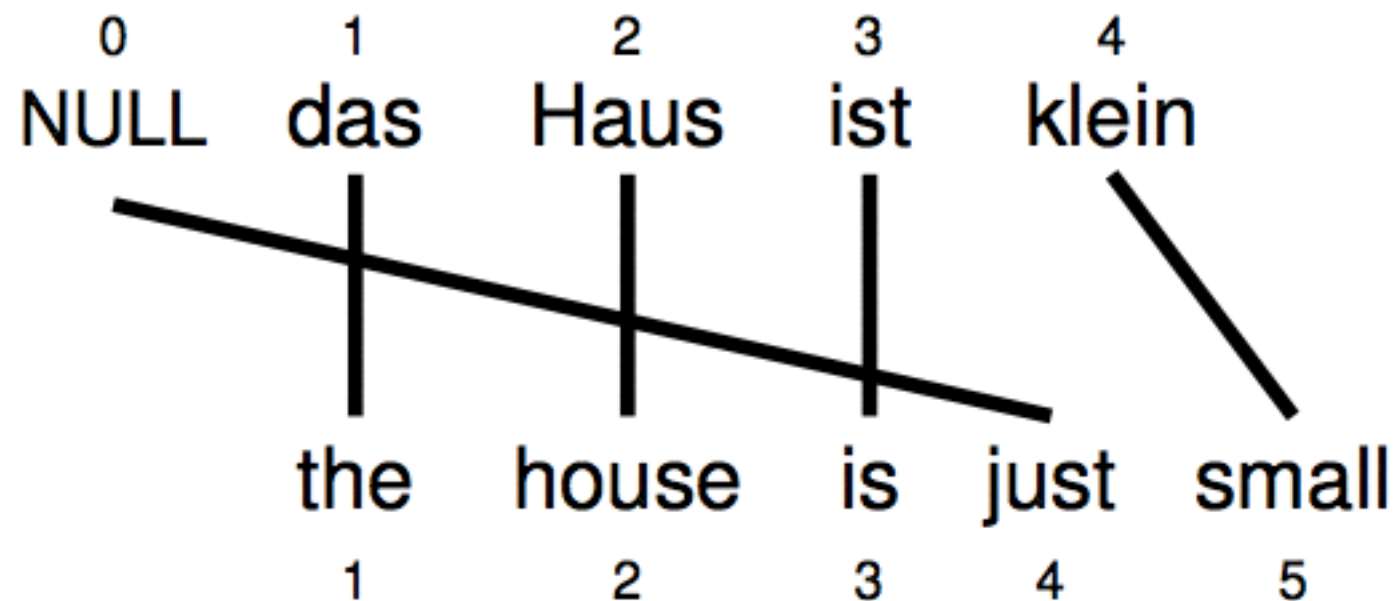


$$\mathbf{a} = (2, 3, 4)^{\top}$$

Word Insertion

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

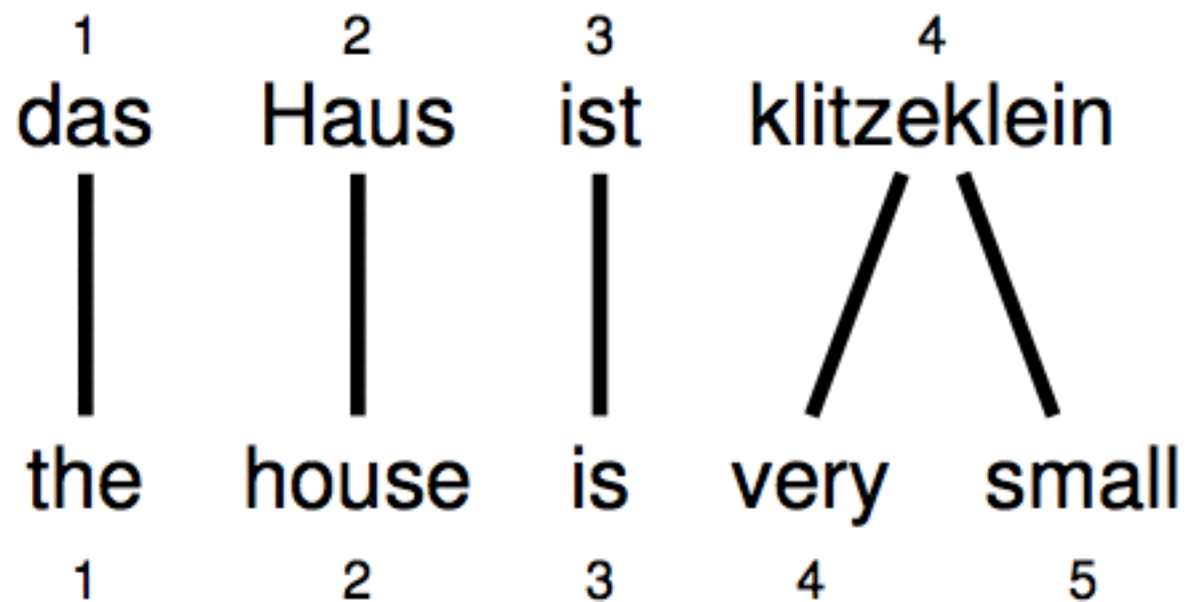
But it must be explained - we typically assume every source sentence contains a NULL token



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

One-to-many Translation

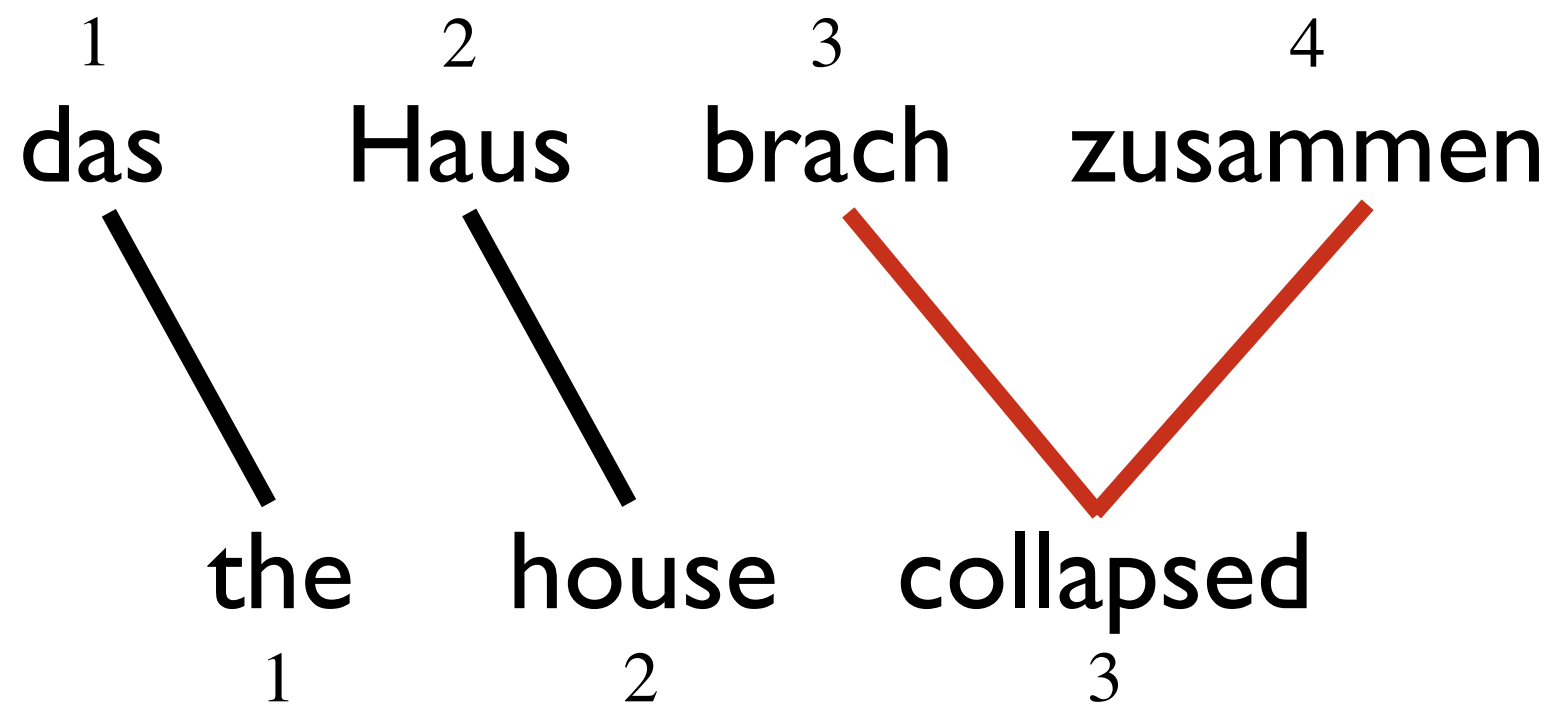
- A source word may translate into **more than one** target word



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

Many-to-one Translation

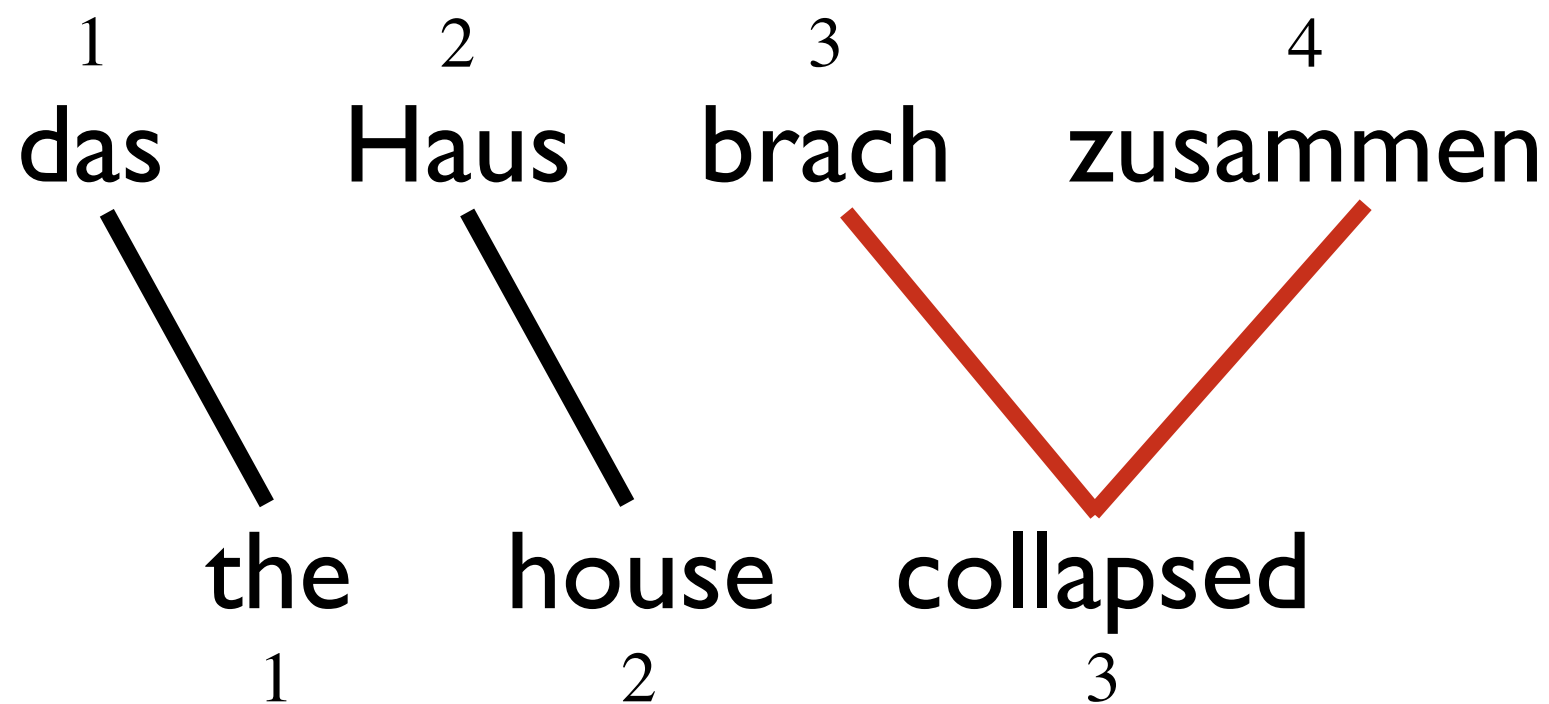
- More than one source word may not translate as a unit in lexical translation



a = ???

Many-to-one Translation

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



$$\mathbf{a} = ???$$

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

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

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

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

Historical Note

IBM Models

Historical Note

IBM Models

The Mathematics of Statistical Machine Translation: Parameter Estimation

Peter F. Brown*
IBM T.J. Watson Research Center

Stephen A. Della Pietra*
IBM T.J. Watson Research Center

Vincent J. Della Pietra*
IBM T.J. Watson Research Center

Robert L. Mercer*
IBM T.J. Watson Research Center

We describe a series of five statistical models of the translation process and give algorithms for estimating the parameters of these models given a set of pairs of sentences that are translations of one another. We define a concept of word-by-word alignment between such pairs of sentences. For any given pair of such sentences each of our models assigns a probability to each of the possible word-by-word alignments. We give an algorithm for seeking the most probable of these alignments. Although the algorithm is suboptimal, the alignment thus obtained accounts well for the word-by-word relationships in the pair of sentences. We have a great deal of data in French and English from the proceedings of the Canadian Parliament. Accordingly, we have restricted our work to these two languages; but we feel that because our algorithms have minimal linguistic content they would work well on other pairs of languages. We also feel, again because of the minimal linguistic content of our algorithms, that it is reasonable to argue that word-by-word alignments are inherent in any sufficiently large bilingual corpus.

1. Introduction

The growing availability of bilingual, machine-readable texts has stimulated interest in methods for extracting linguistically valuable information from such texts. For ex-

Historical Note

IBM Models

The Mathematics of Statistical Machine Translation: Parameter Estimation

Peter F. Brown*
IBM T.J. Watson Research Center

Stephen A. Della Pietra*
IBM T.J. Watson Research Center

Vincent J. Della Pietra*
IBM T.J. Watson Research Center

Robert L. Mercer*
IBM T.J. Watson Research Center

We describe a series of five statistical models of the translation process and give algorithms for estimating the parameters of these models given a set of pairs of sentences that are translations of one another. We define a concept of word-by-word alignment between such pairs of sentences. For any given pair of such sentences each of our models assigns a probability to each of the possible word-by-word alignments. We give an algorithm for seeking the most probable of these alignments. Although the algorithm is suboptimal, the alignment thus obtained accounts well for the word-by-word relationships in the pair of sentences. We have a great deal of data in French and English from the proceedings of the Canadian Parliament. Accordingly, we have restricted our work to these two languages; but we feel that because our algorithms have minimal linguistic content they would work well on other pairs of languages. We also feel, again because of the minimal linguistic content of our algorithms, that it is reasonable to argue that word-by-word alignments are inherent in any sufficiently large bilingual corpus.

1. Introduction

The growing availability of bilingual, machine-readable texts has stimulated interest in methods for extracting linguistically valuable information from such texts. For ex-



Fred Jelinek
(1932-2010)

Historical Note

IBM Models

Some of us started to wonder in the mid 1980s whether our [speech recognition] methods could be successfully applied to new fields. Bob Mercer and I spent many of our after-lunch “periphery” walks discussing possible candidates. We soon came up with two: machine translation and stock market modeling



Fred Jelinek
(1932-2010)

minimal linguistic content of our algorithms, that it is reasonable to argue that word-by-word alignments are inherent in any sufficiently large bilingual corpus.

1. Introduction

The growing availability of bilingual, machine-readable texts has stimulated interest in methods for extracting linguistically valuable information from such texts. For ex-

Historical Note

IBM Models

The Mathematics of Statistical Machine Translation: Parameter Estimation

Peter F. Brown*
IBM T.J. Watson Research Center

Stephen A. Della Pietra*
IBM T.J. Watson Research Center

Vincent J. Della Pietra*
IBM T.J. Watson Research Center

Robert L. Mercer*
IBM T.J. Watson Research Center

We describe a series of five statistical models of the translation process and give algorithms for estimating the parameters of these models given a set of pairs of sentences that are translations of one another. We define a concept of word-by-word alignment between such pairs of sentences. For any given pair of such sentences each of our models assigns a probability to each of the possible word-by-word alignments. We give an algorithm for seeking the most probable of these alignments. Although the algorithm is suboptimal, the alignment thus obtained accounts well for the word-by-word relationships in the pair of sentences. We have a great deal of data in French and English from the proceedings of the Canadian Parliament. Accordingly, we have restricted our work to these two languages; but we feel that because our algorithms have minimal linguistic content they would work well on other pairs of languages. We also feel, again because of the minimal linguistic content of our algorithms, that it is reasonable to argue that word-by-word alignments are inherent in any sufficiently large bilingual corpus.

1. Introduction

The growing availability of bilingual, machine-readable texts has stimulated interest in methods for extracting linguistically valuable information from such texts. For ex-



Fred Jelinek
(1932-2010)

Historical Note

IBM Models

“The validity of a 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, p. 30ff and references therein). The crude force of computers is not science. The paper is simply beyond the scope of COLING.”



Fred Jelinek
(1932-2010)

Historical Note

IBM Models

The Mathematics of Statistical Machine Translation: Parameter Estimation

Peter F. Brown*
IBM T.J. Watson Research Center

Stephen A. Della Pietra*
IBM T.J. Watson Research Center

Vincent J. Della Pietra*
IBM T.J. Watson Research Center

Robert L. Mercer*
IBM T.J. Watson Research Center

We describe a series of five statistical models of the translation process and give algorithms for estimating the parameters of these models given a set of pairs of sentences that are translations of one another. We define a concept of word-by-word alignment between such pairs of sentences. For any given pair of such sentences each of our models assigns a probability to each of the possible word-by-word alignments. We give an algorithm for seeking the most probable of these alignments. Although the algorithm is suboptimal, the alignment thus obtained accounts well for the word-by-word relationships in the pair of sentences. We have a great deal of data in French and English from the proceedings of the Canadian Parliament. Accordingly, we have restricted our work to these two languages; but we feel that because our algorithms have minimal linguistic content they would work well on other pairs of languages. We also feel, again because of the minimal linguistic content of our algorithms, that it is reasonable to argue that word-by-word alignments are inherent in any sufficiently large bilingual corpus.

1. Introduction

The growing availability of bilingual, machine-readable texts has stimulated interest in methods for extracting linguistically valuable information from such texts. For ex-



Fred Jelinek
(1932-2010)

Historical Note

IBM Models

Renaissance



The Mathematics of Statistical Machine Translation: Parameter Estimation

Peter F. Brown*
IBM T.J. Watson Research Center

Stephen A. Della Pietra*
IBM T.J. Watson Research Center

Vincent J. Della Pietra*
IBM T.J. Watson Research Center

Robert L. Mercer*
IBM T.J. Watson Research Center

We describe a series of five statistical models of the translation process and give algorithms for estimating the parameters of these models given a set of pairs of sentences that are translations of one another. We define a concept of word-by-word alignment between such pairs of sentences. For any given pair of such sentences each of our models assigns a probability to each of the possible word-by-word alignments. We give an algorithm for seeking the most probable of these alignments. Although the algorithm is suboptimal, the alignment thus obtained accounts well for the word-by-word relationships in the pair of sentences. We have a great deal of data in French and English from the proceedings of the Canadian Parliament. Accordingly, we have restricted our work to these two languages; but we feel that because our algorithms have minimal linguistic content they would work well on other pairs of languages. We also feel, again because of the minimal linguistic content of our algorithms, that it is reasonable to argue that word-by-word alignments are inherent in any sufficiently large bilingual corpus.

1. Introduction

The growing availability of bilingual, machine-readable texts has stimulated interest in methods for extracting linguistically valuable information from such texts. For ex-



Fred Jelinek
(1932-2010)

Historical Note

IBM Models

Renaissance



The Mathematics of Statistical Machine Translation: Parameter Estimation

Peter F. Brown*
IBM T.J. Watson Research Center

Stephen A. Della Pietra*
IBM T.J. Watson Research Center

Vincent J. Della Pietra*
IBM T.J. Watson Research Center

Robert L. Mercer*
IBM T.J. Watson Research Center

We describe a series of five statistical models of the translation process and give algorithms for estimating the parameters of these models given a set of pairs of sentences that are translations of one another. We define a concept of word-by-word alignment between such pairs of sentences. For any given pair of such sentences each of our models assigns a probability to each of the possible word-by-word alignments. We give an algorithm for seeking the most probable of these alignments. Although the algorithm is suboptimal, the alignment thus obtained accounts well for the word-by-word relationships in the pair of sentences. We have a great deal of data in French and English from the proceedings of the Canadian Parliament. Accordingly, we have restricted our work to these two languages; but we feel that because our algorithms have minimal linguistic content they would work well on other pairs of languages. We also feel, again because of the minimal linguistic content of our algorithms, that it is reasonable to argue that word-by-word alignments are inherent in any sufficiently large bilingual corpus.

1. Introduction

The growing availability of bilingual, machine-readable texts has stimulated interest in methods for extracting linguistically valuable information from such texts. For ex-



Fred Jelinek
(1932-2010)



**The Center For Language
and Speech Processing**
at the Johns Hopkins University

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

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

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

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

IBM Model I

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

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

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

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

IBM Model I

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

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

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

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

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

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

Marginal probability

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

$$p(e_i \mid \mathbf{f}, m) = \sum_{a_i=0}^n \frac{1}{1+n} p(e_i \mid f_{a_i})$$

Recall our independence assumption: all alignment decisions are independent of each other, and given the alignments then all translation decisions are independent of each other, so **all translation decisions are independent of each other.**

Marginal probability

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

$$p(e_i \mid \mathbf{f}, m) = \sum_{a_i=0}^n \frac{1}{1+n} p(e_i \mid f_{a_i})$$

Recall our independence assumption: all alignment decisions are independent of each other, and given the alignments then all translation decisions are independent of each other, so **all translation decisions are independent of each other.**

$$p(a, b, c, d) = p(a)p(b)p(c)p(d)$$

Marginal probability

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

$$p(e_i \mid \mathbf{f}, m) = \sum_{a_i=0}^n \frac{1}{1+n} p(e_i \mid f_{a_i})$$

Recall our independence assumption: all alignment decisions are independent of each other, and given the alignments then all translation decisions are independent of each other, so **all translation decisions are independent of each other.**

Marginal probability

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

$$p(e_i \mid \mathbf{f}, m) = \sum_{a_i=0}^n \frac{1}{1+n} p(e_i \mid f_{a_i})$$

Recall our independence assumption: all alignment decisions are independent of each other, and given the alignments then all translation decisions are independent of each other, so **all translation decisions are independent of each other.**

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

Marginal probability

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

$$p(e_i \mid \mathbf{f}, m) = \sum_{a_i=0}^n \frac{1}{1+n} p(e_i \mid f_{a_i})$$

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

Marginal probability

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

$$p(e_i \mid \mathbf{f}, m) = \sum_{a_i=0}^n \frac{1}{1+n} p(e_i \mid f_{a_i})$$

$$\begin{aligned} p(\mathbf{e} \mid \mathbf{f}, m) &= \prod_{i=1}^m p(e_i \mid \mathbf{f}, m) \\ &= \prod_{i=1}^m \sum_{a_i=0}^n \frac{1}{1+n} p(e_i \mid f_{a_i}) \end{aligned}$$

Marginal probability

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

$$p(e_i \mid \mathbf{f}, m) = \sum_{a_i=0}^n \frac{1}{1+n} p(e_i \mid f_{a_i})$$

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

$$= \prod_{i=1}^m \sum_{a_i=0}^n \frac{1}{1+n} p(e_i \mid f_{a_i})$$

$$= \frac{1}{(1+n)^m} \prod_{i=1}^m \sum_{a_i=0}^n p(e_i \mid f_{a_i})$$

Example

0	1	2	3	4
NULL	das	Haus	ist	klein

1	2	3	4
---	---	---	---

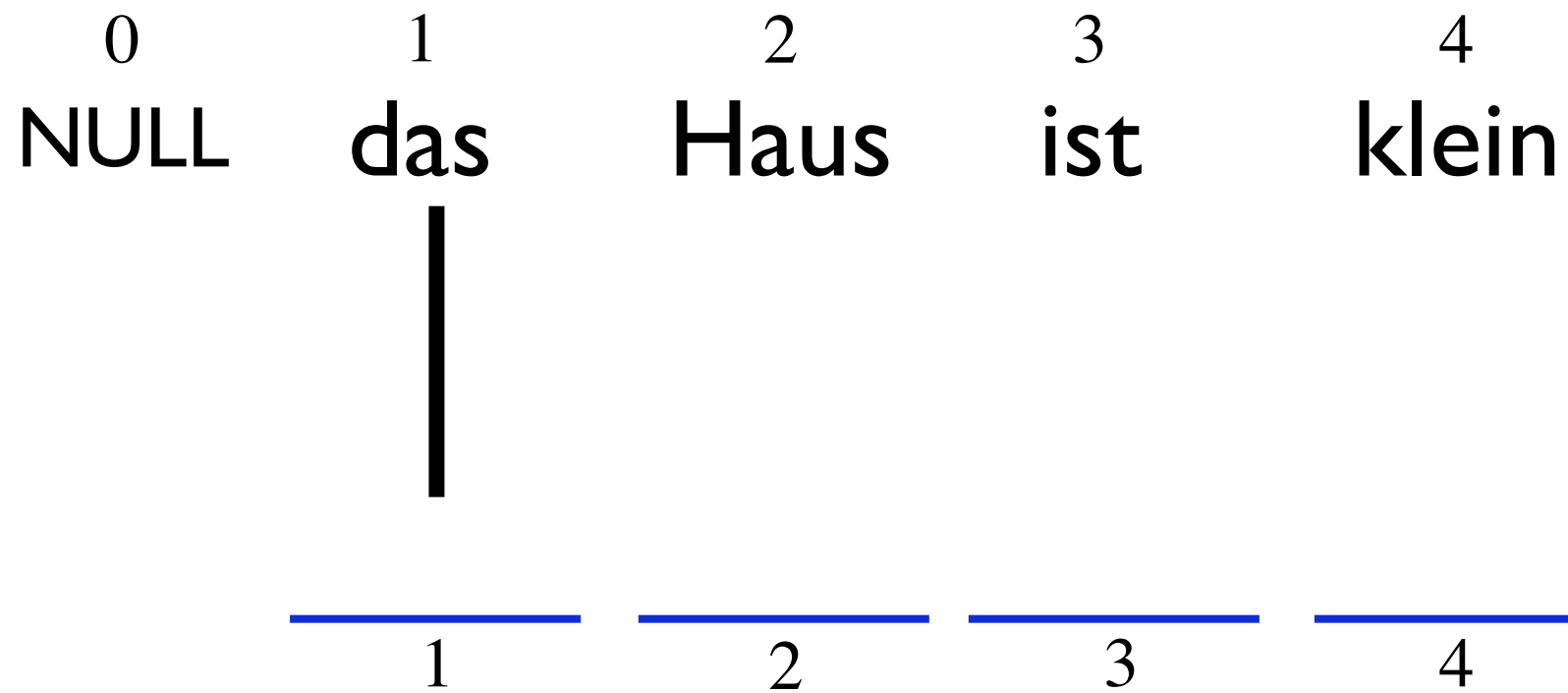
Start with a foreign sentence and a target length.

Example

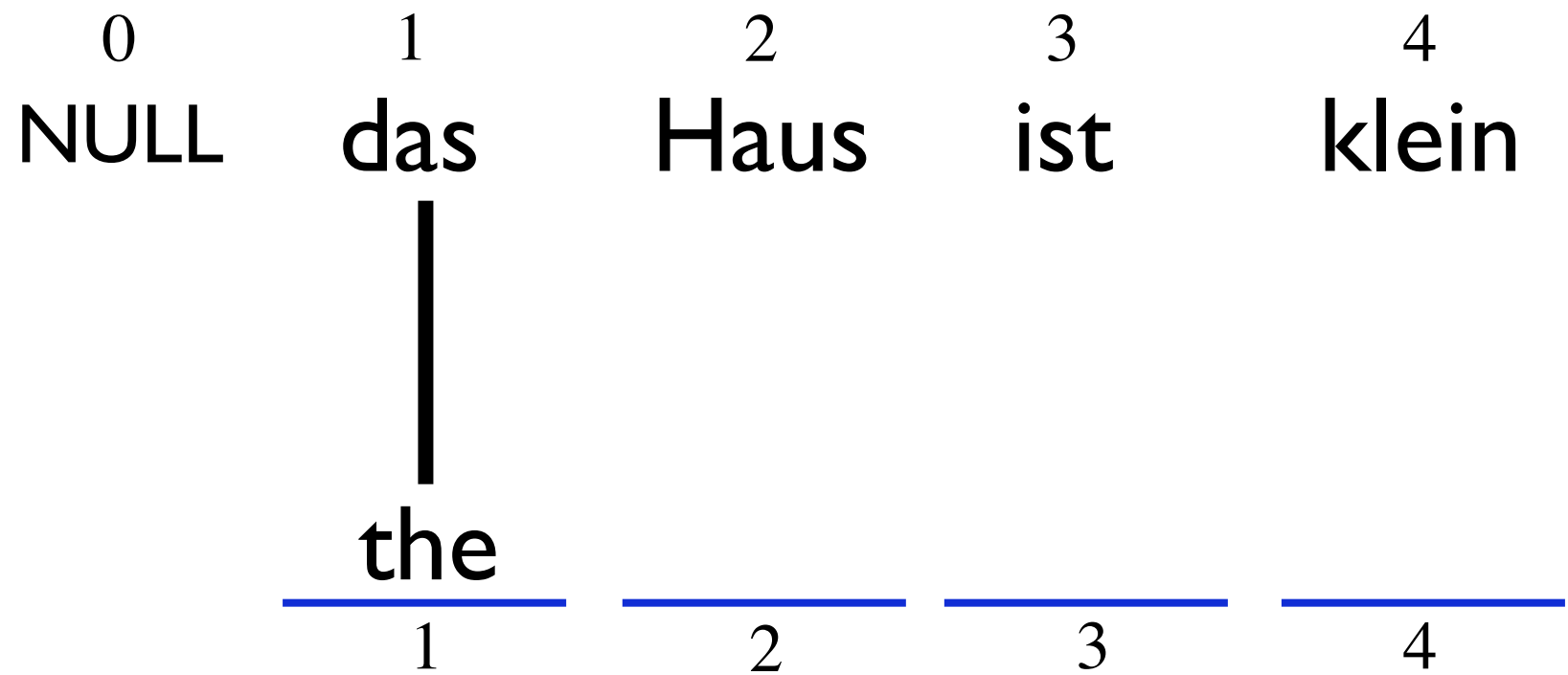
0	1	2	3	4
NULL	das	Haus	ist	klein

1	2	3	4
---	---	---	---

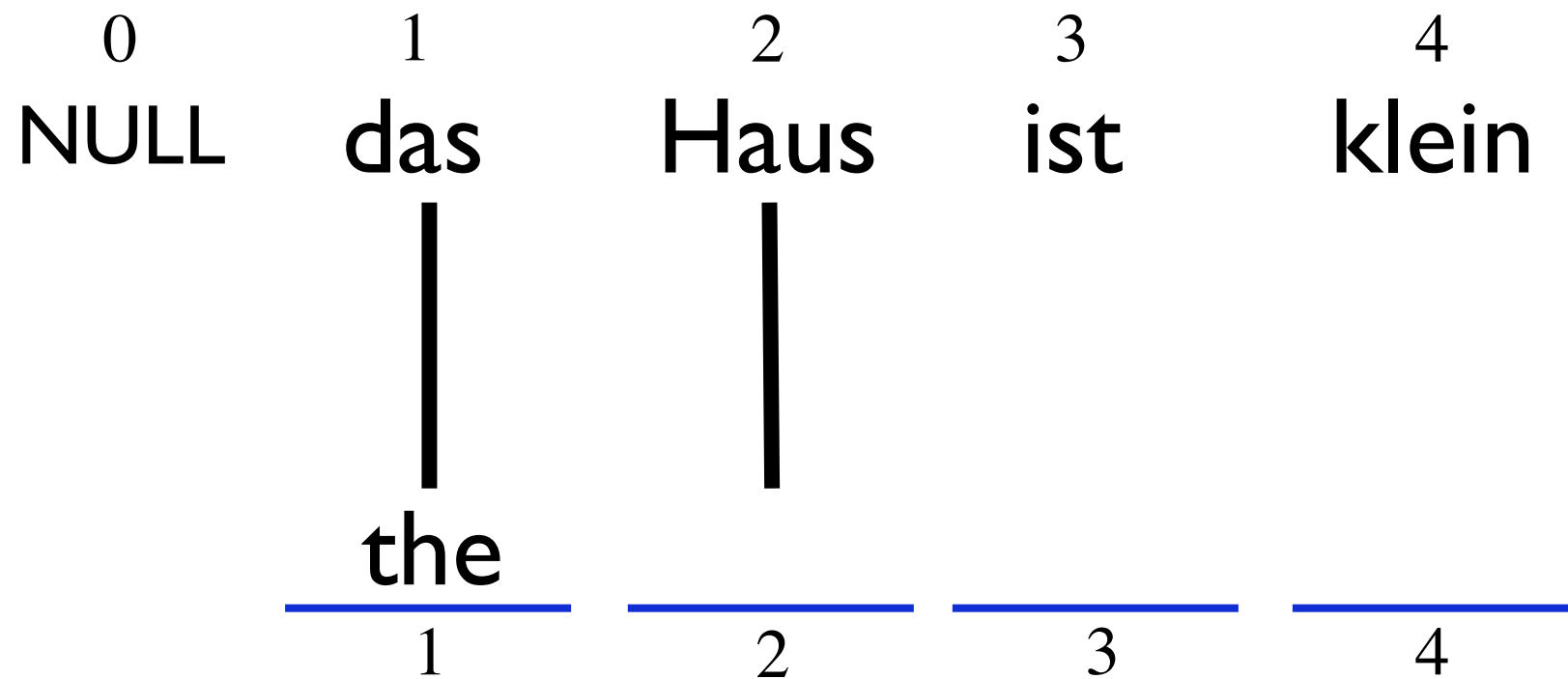
Example



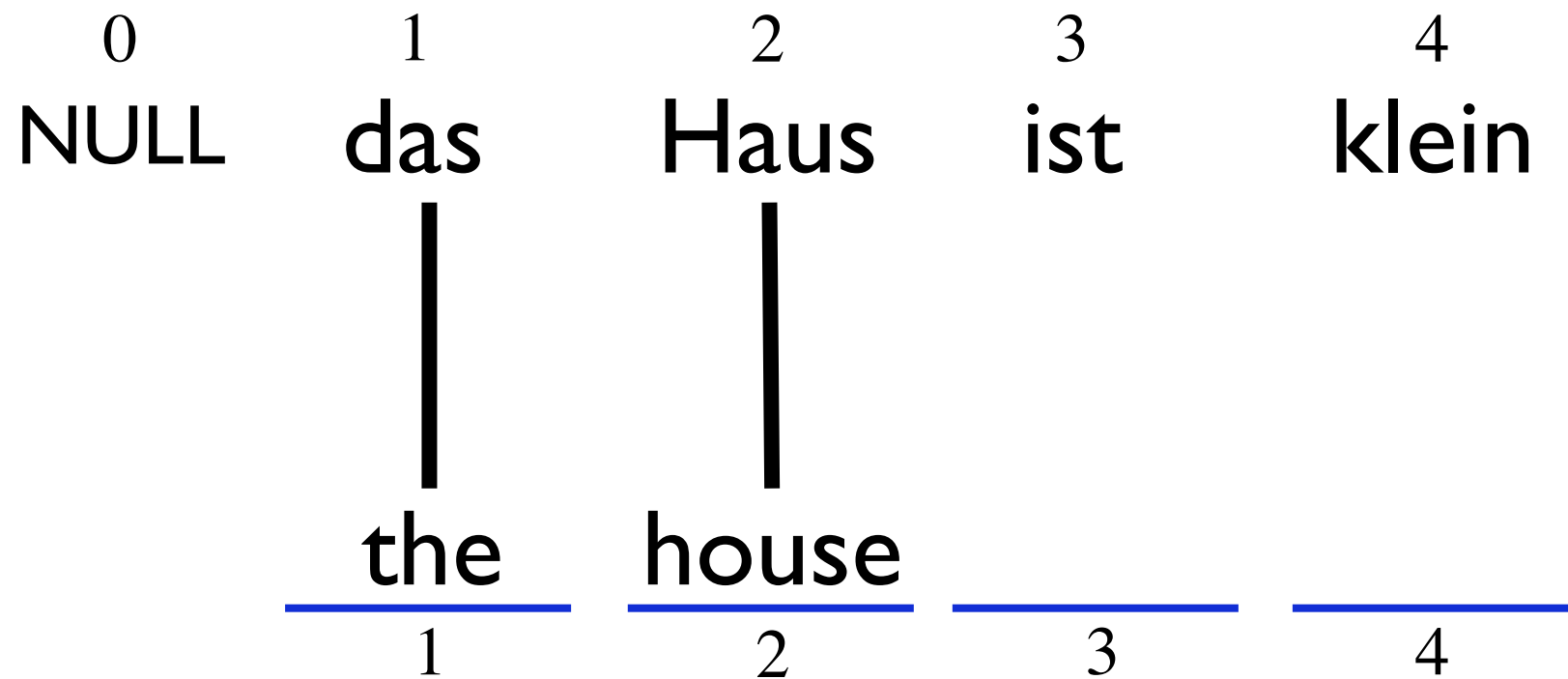
Example



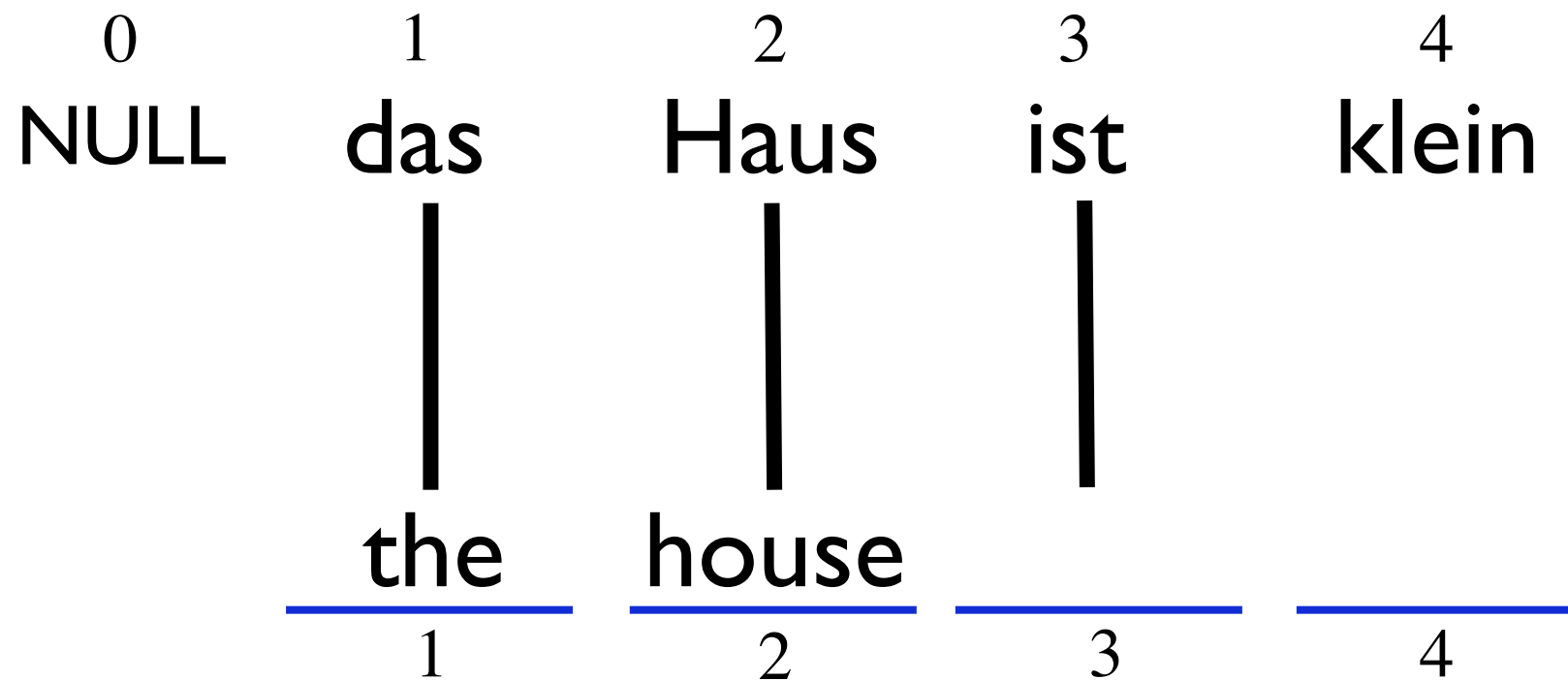
Example



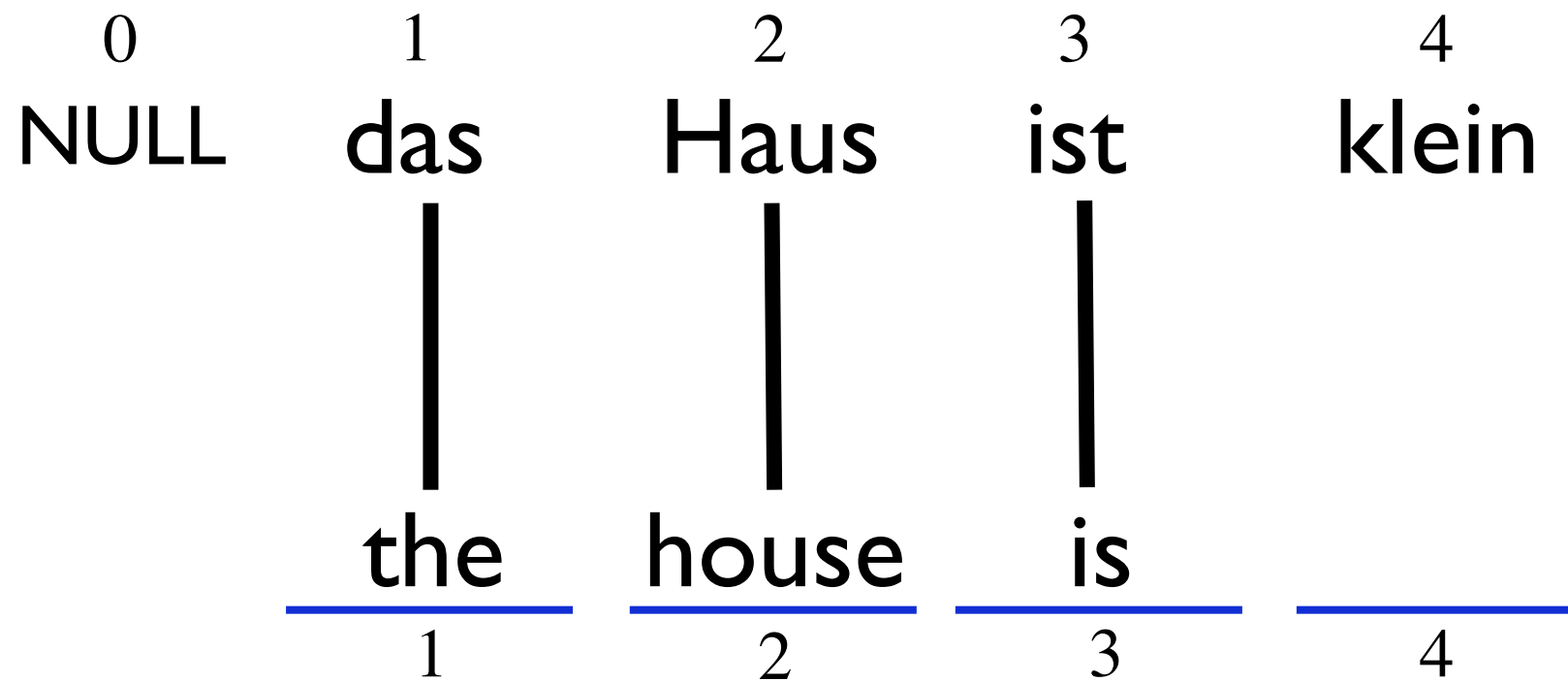
Example



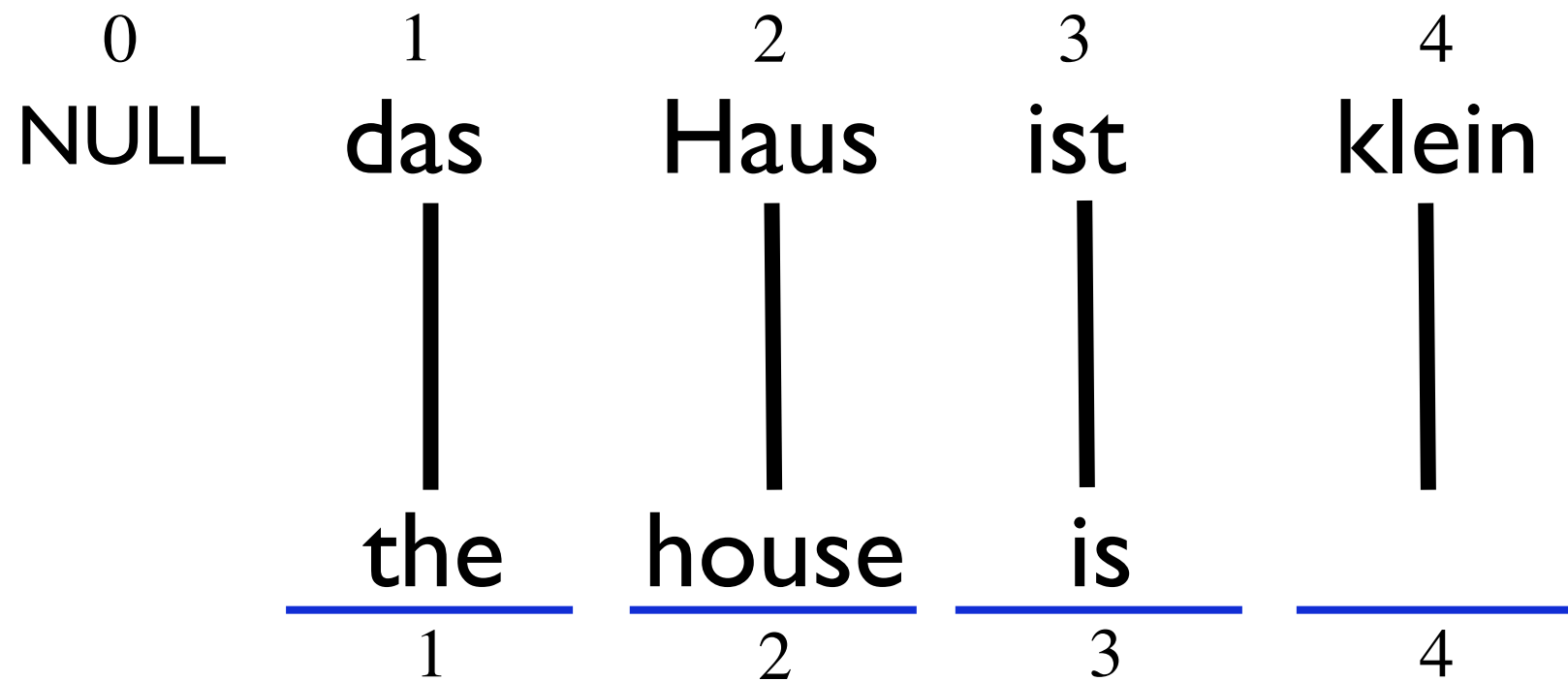
Example



Example



Example



Example

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

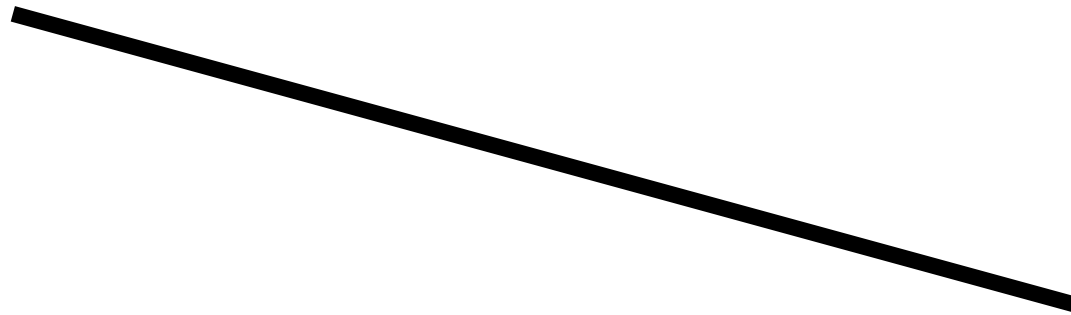
Example

0	1	2	3	4
NULL	das	Haus	ist	klein

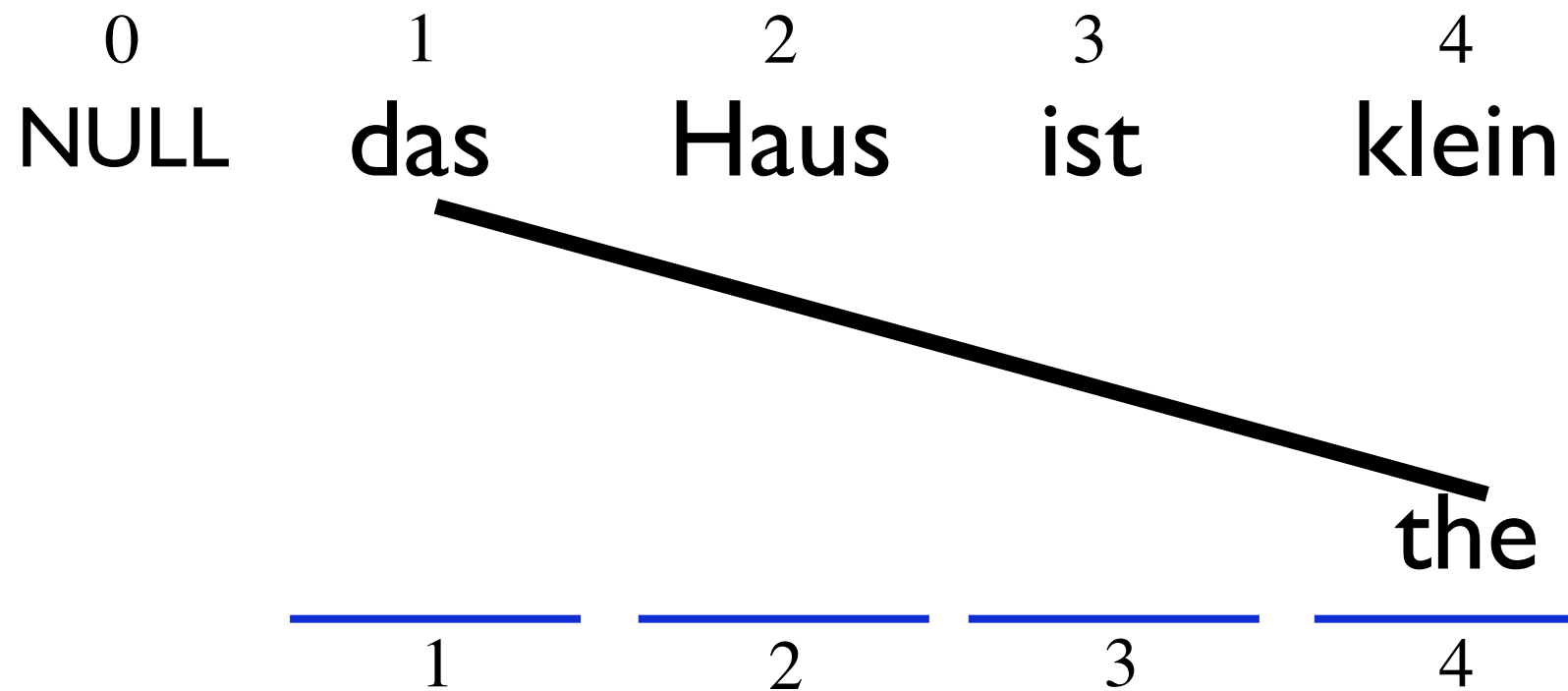
1	2	3	4
---	---	---	---

Example

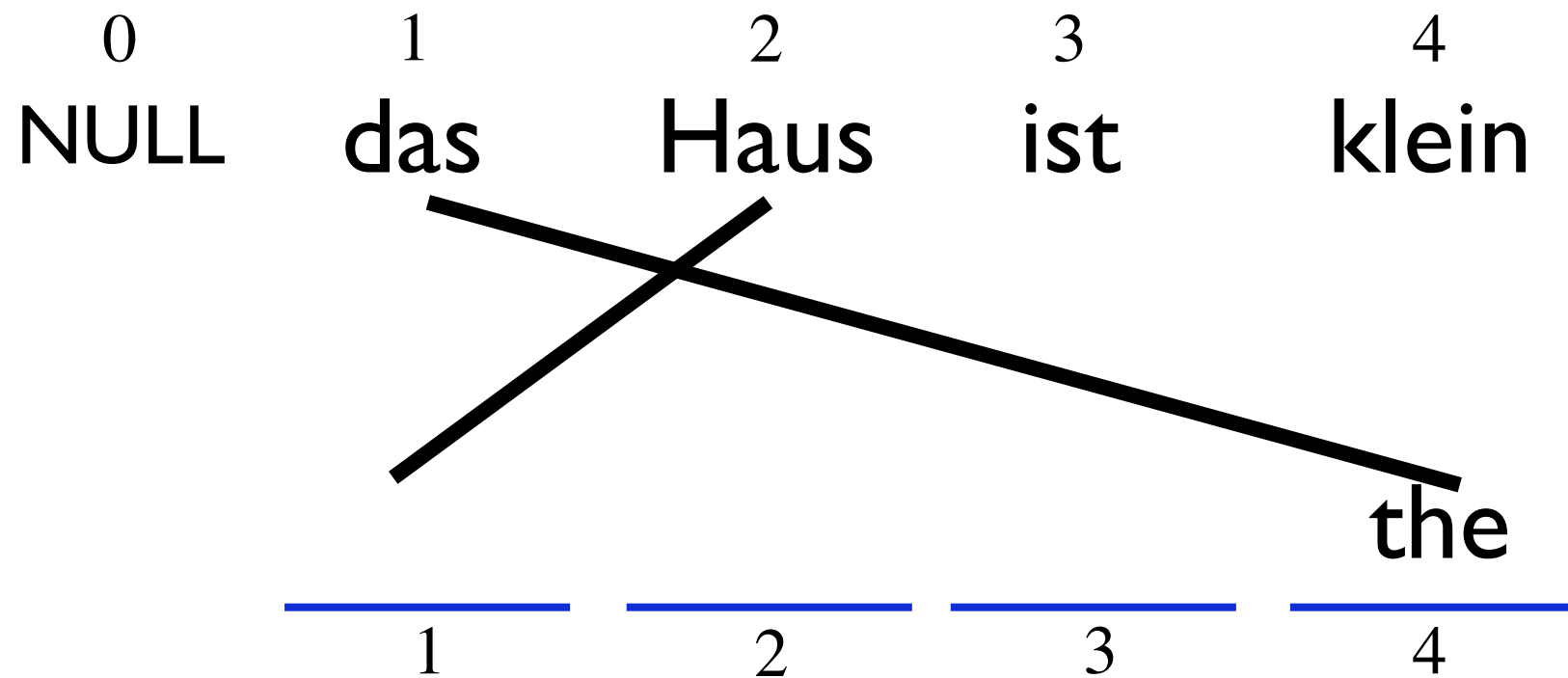
0	1	2	3	4
NULL	das	Haus	ist	klein



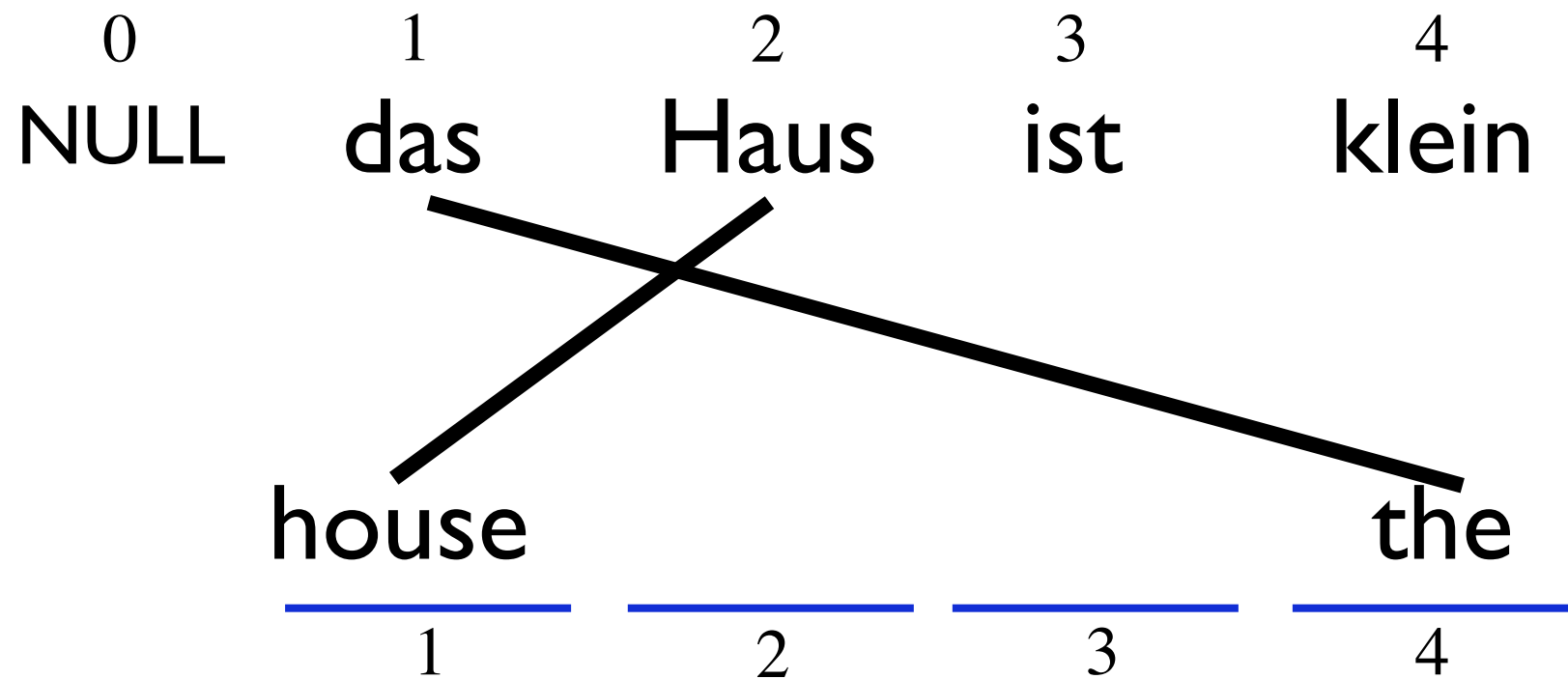
Example



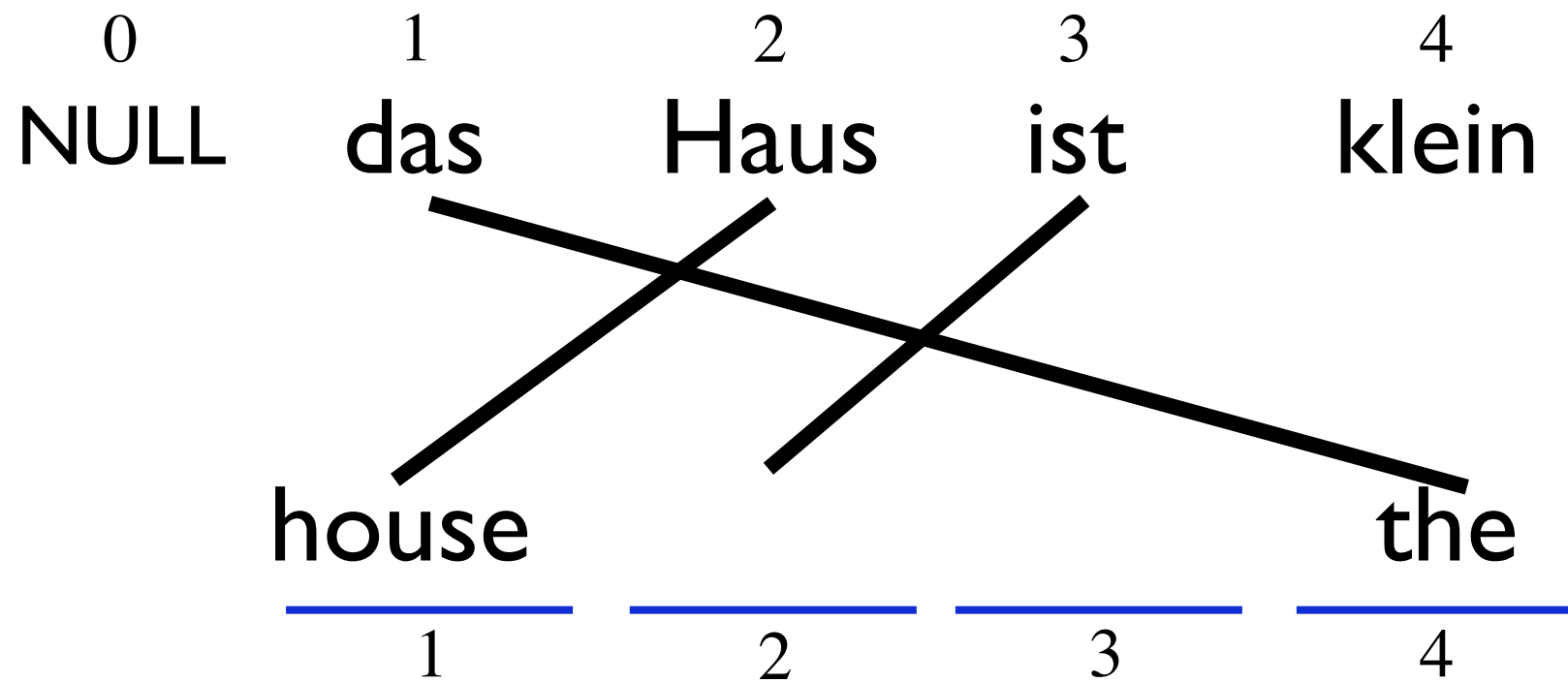
Example



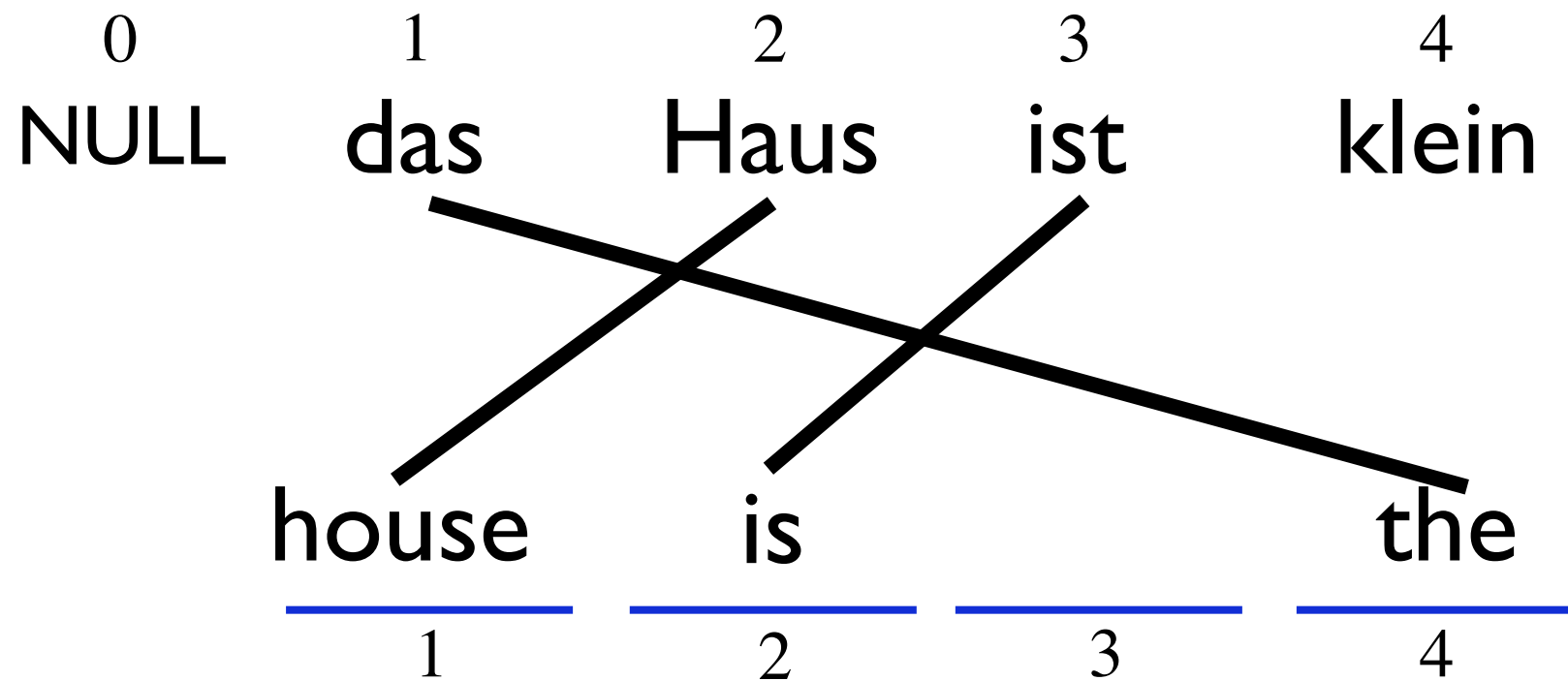
Example



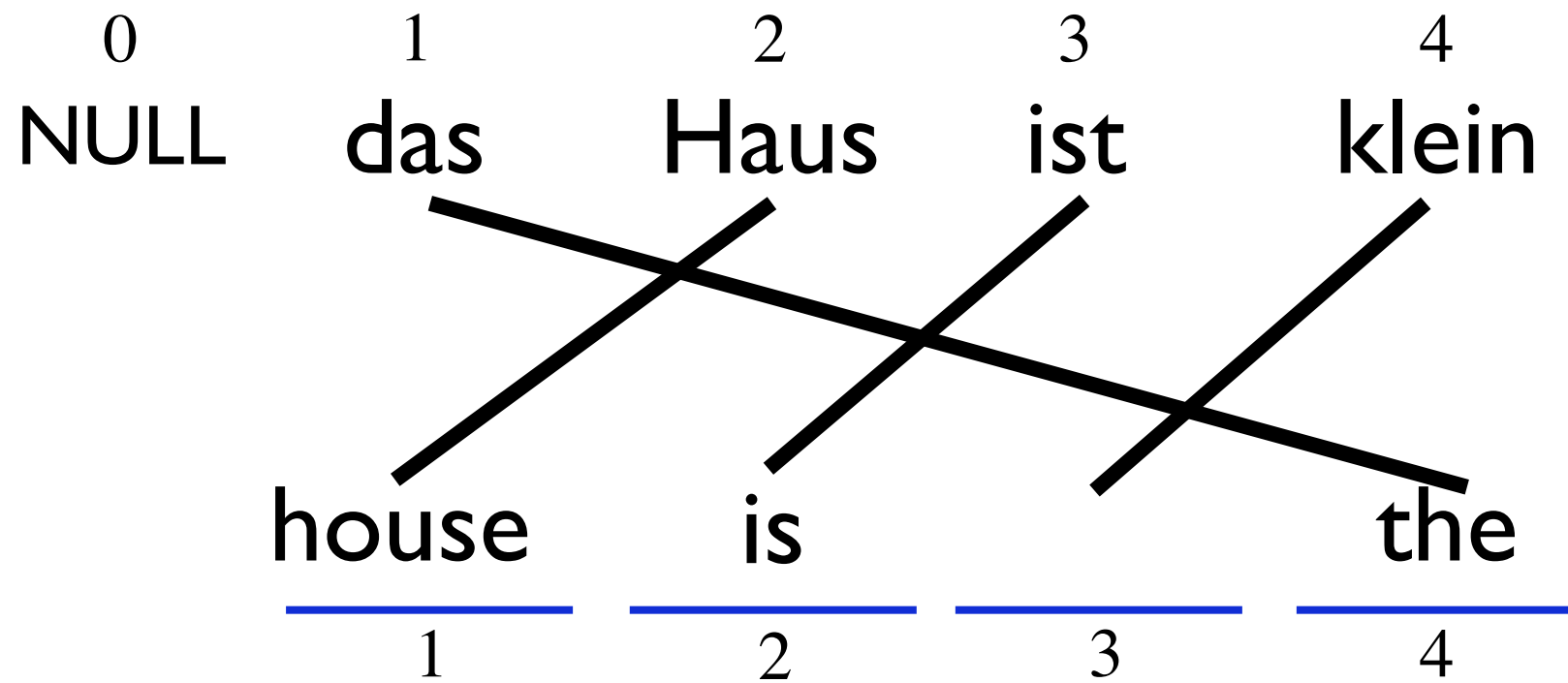
Example



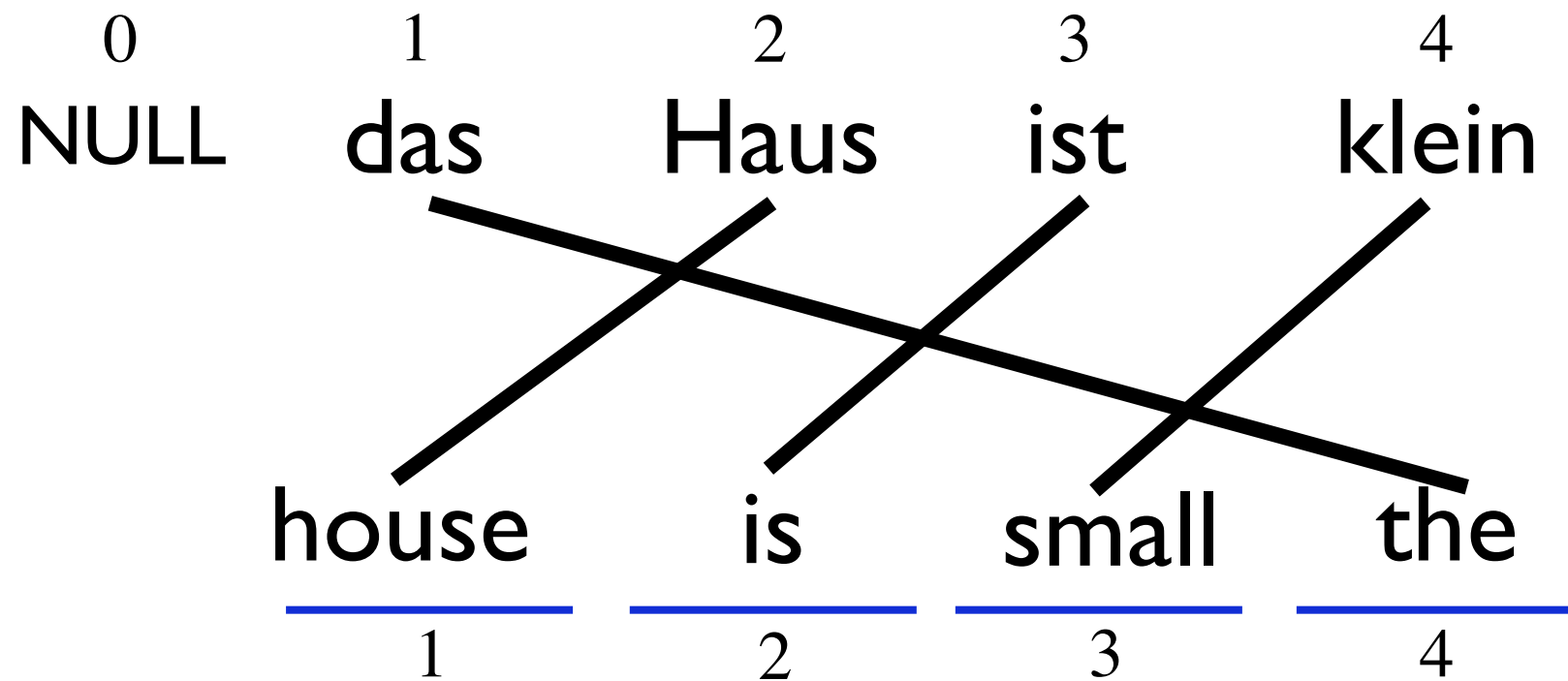
Example



Example



Example



Finding the Viterbi Alignment

$$\mathbf{a}^* = \arg \max_{\mathbf{a} \in [0,1,\dots,n]^m} p(\mathbf{a} \mid \mathbf{e}, \mathbf{f})$$

Finding the Viterbi Alignment

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

Finding the Viterbi Alignment

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

Finding the Viterbi Alignment

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

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

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 #2

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

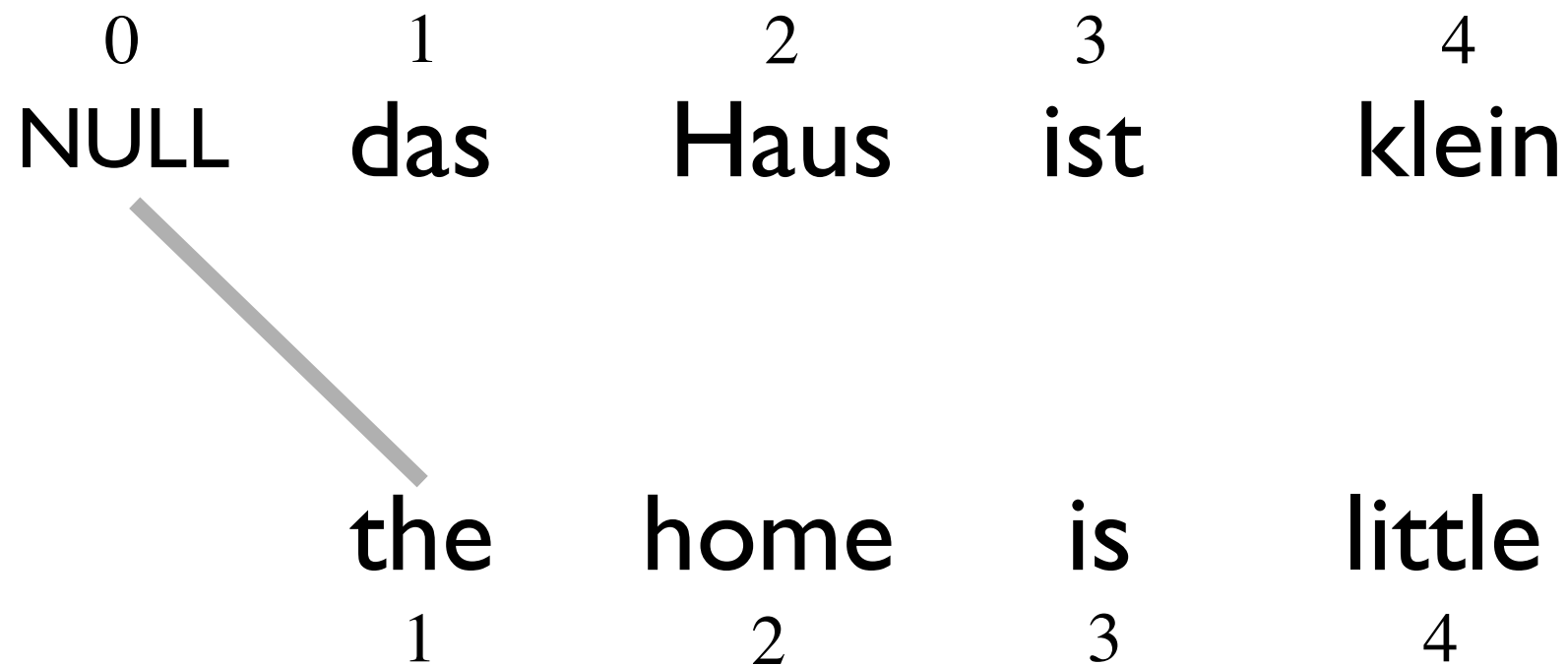


Finding the Viterbi Alignment

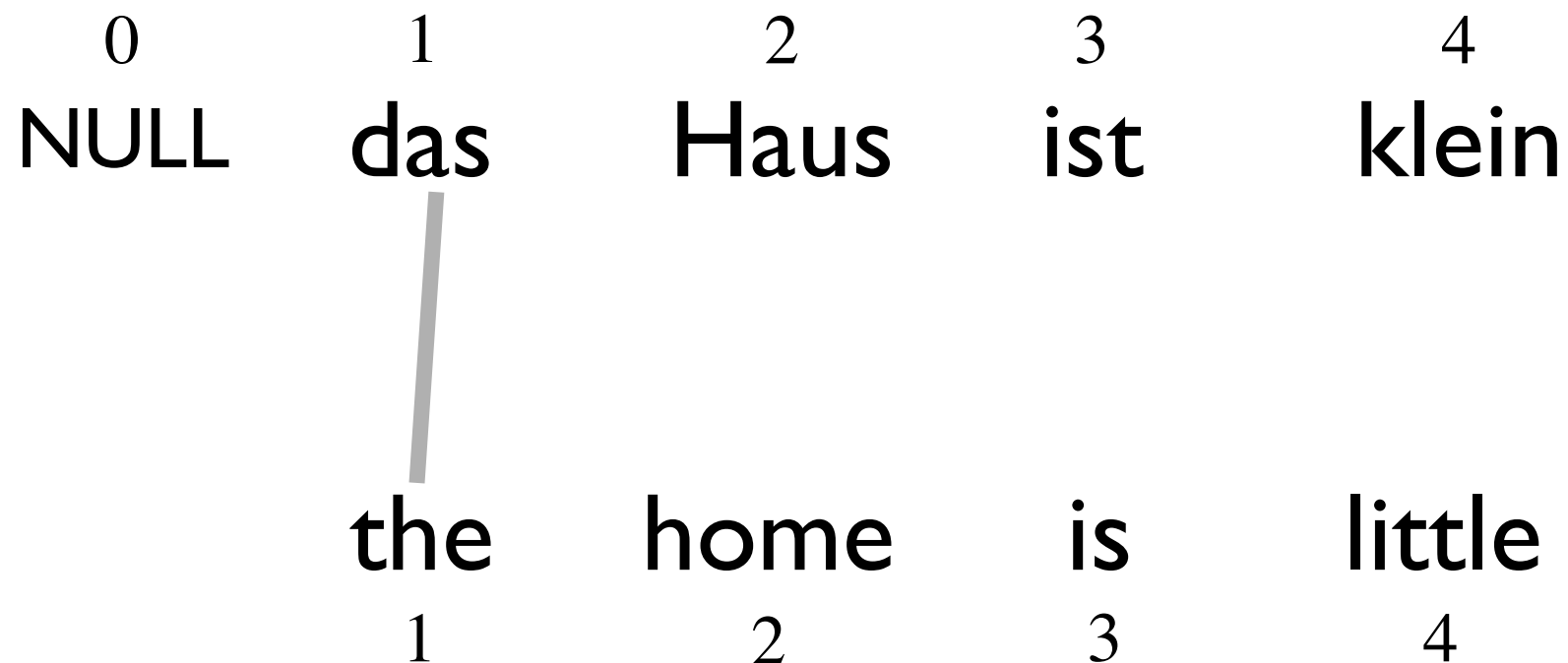
0	1	2	3	4
NULL	das	Haus	ist	klein

the	home	is	little
1	2	3	4

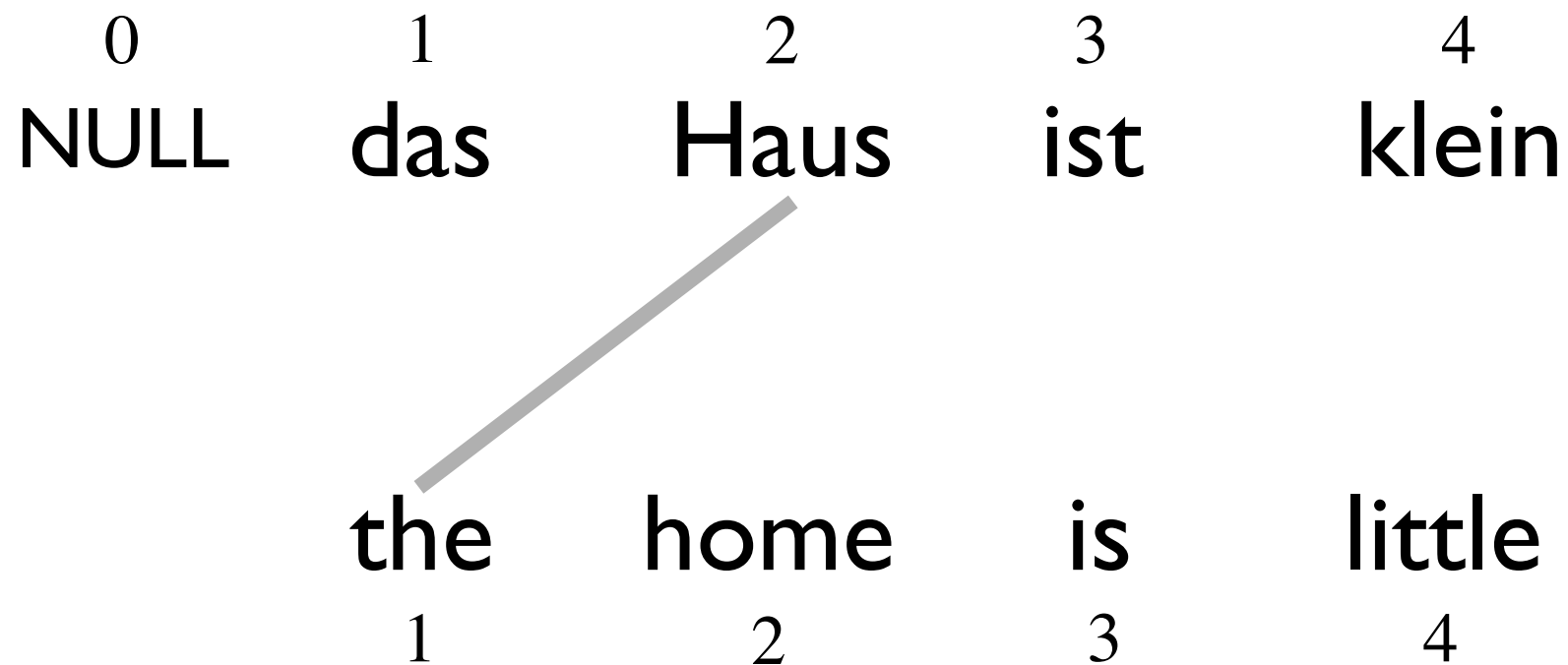
Finding the Viterbi Alignment



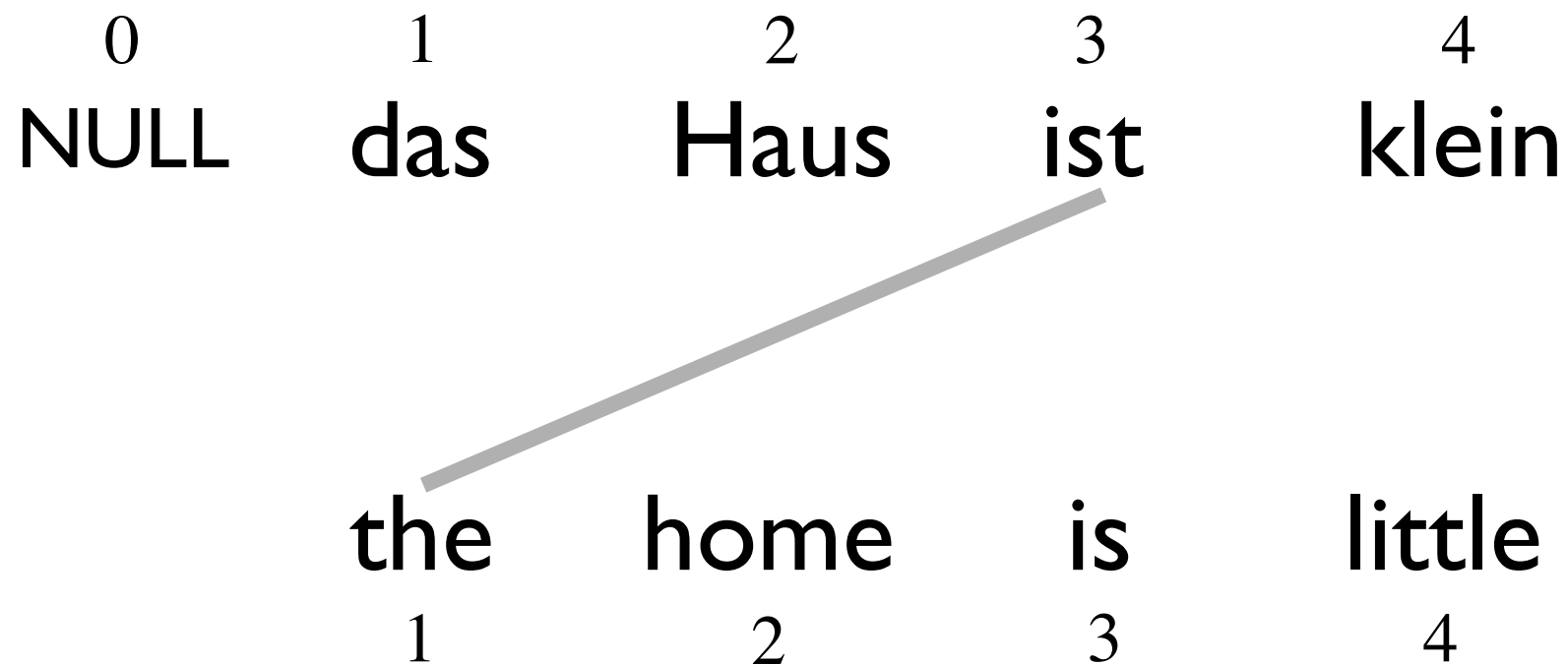
Finding the Viterbi Alignment



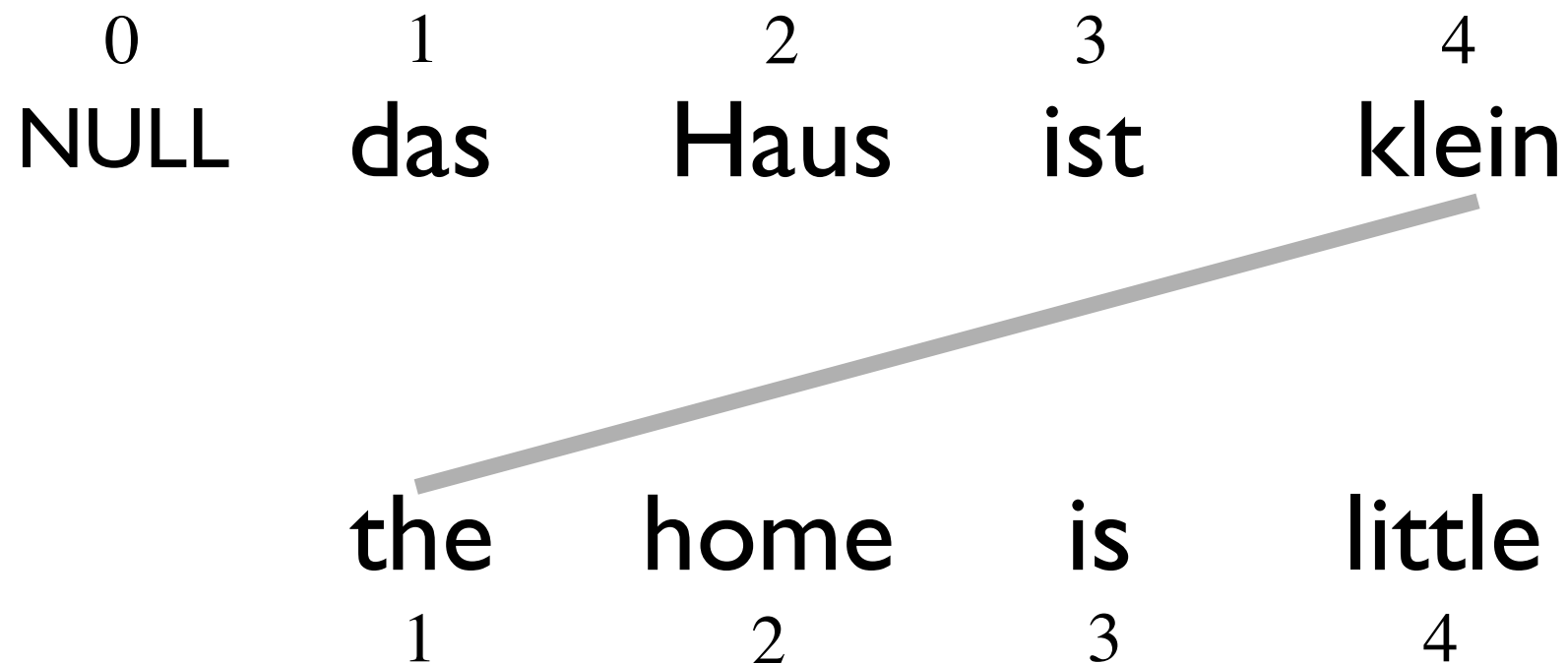
Finding the Viterbi Alignment



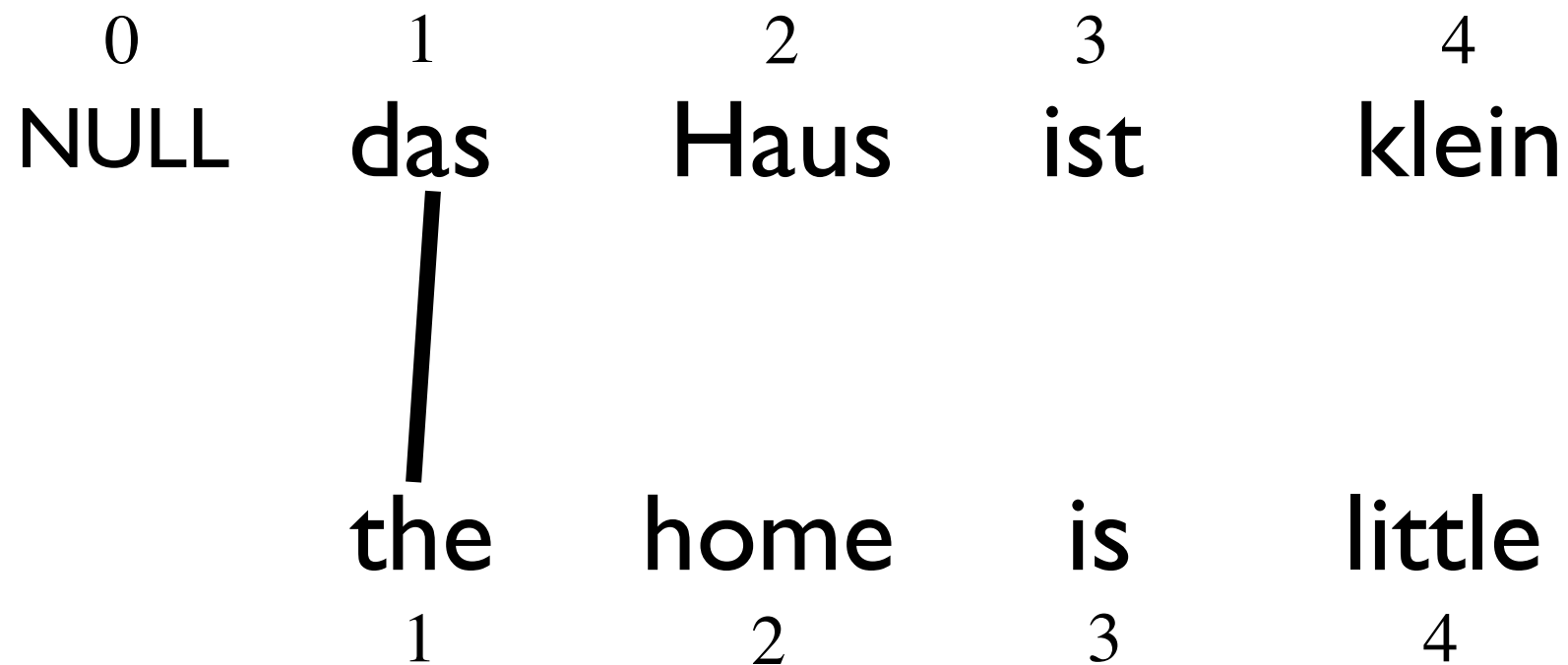
Finding the Viterbi Alignment



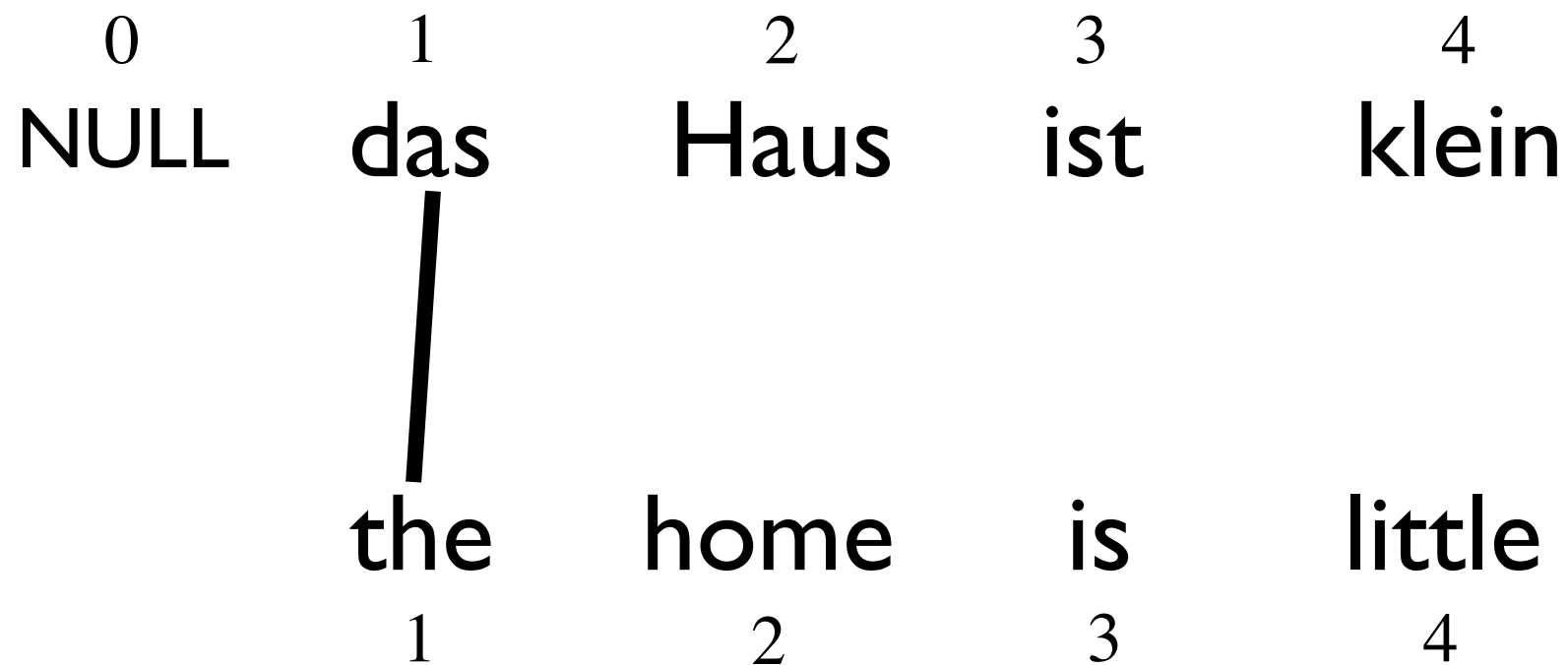
Finding the Viterbi Alignment



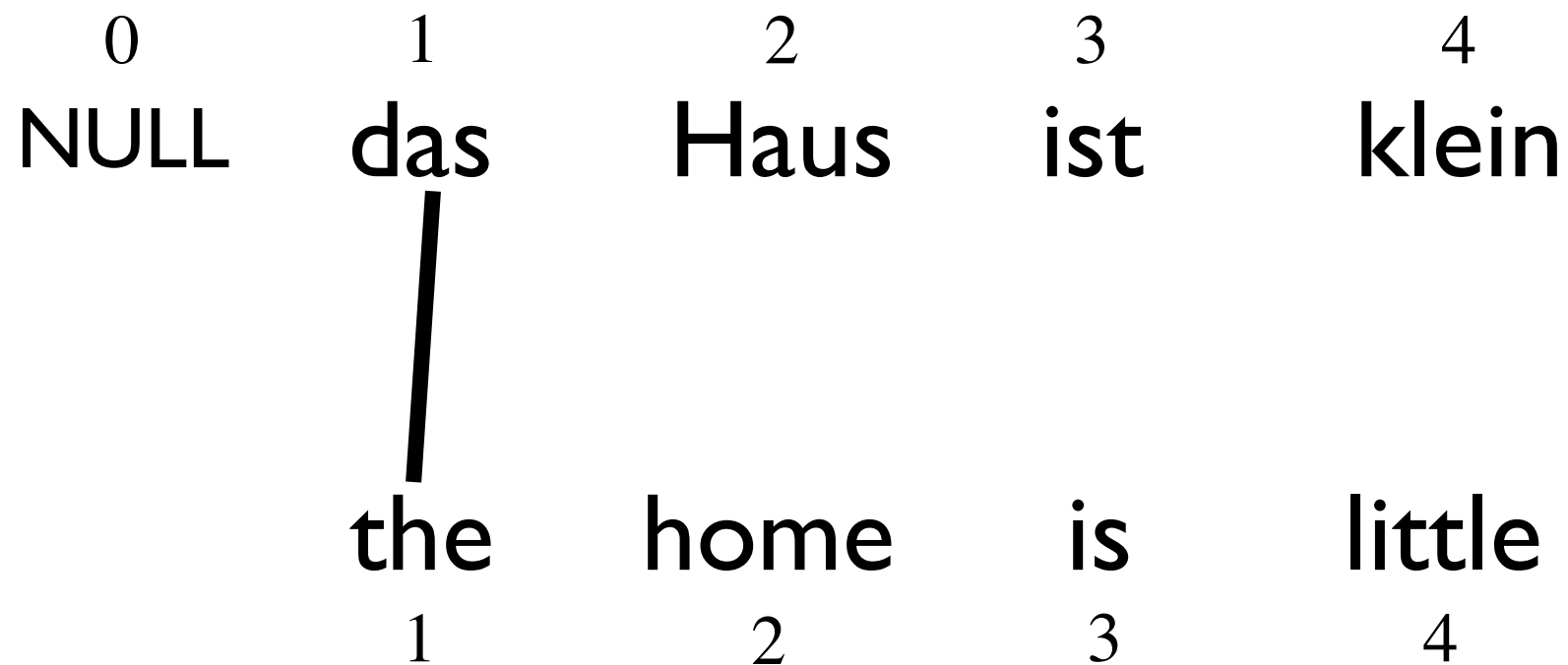
Finding the Viterbi Alignment



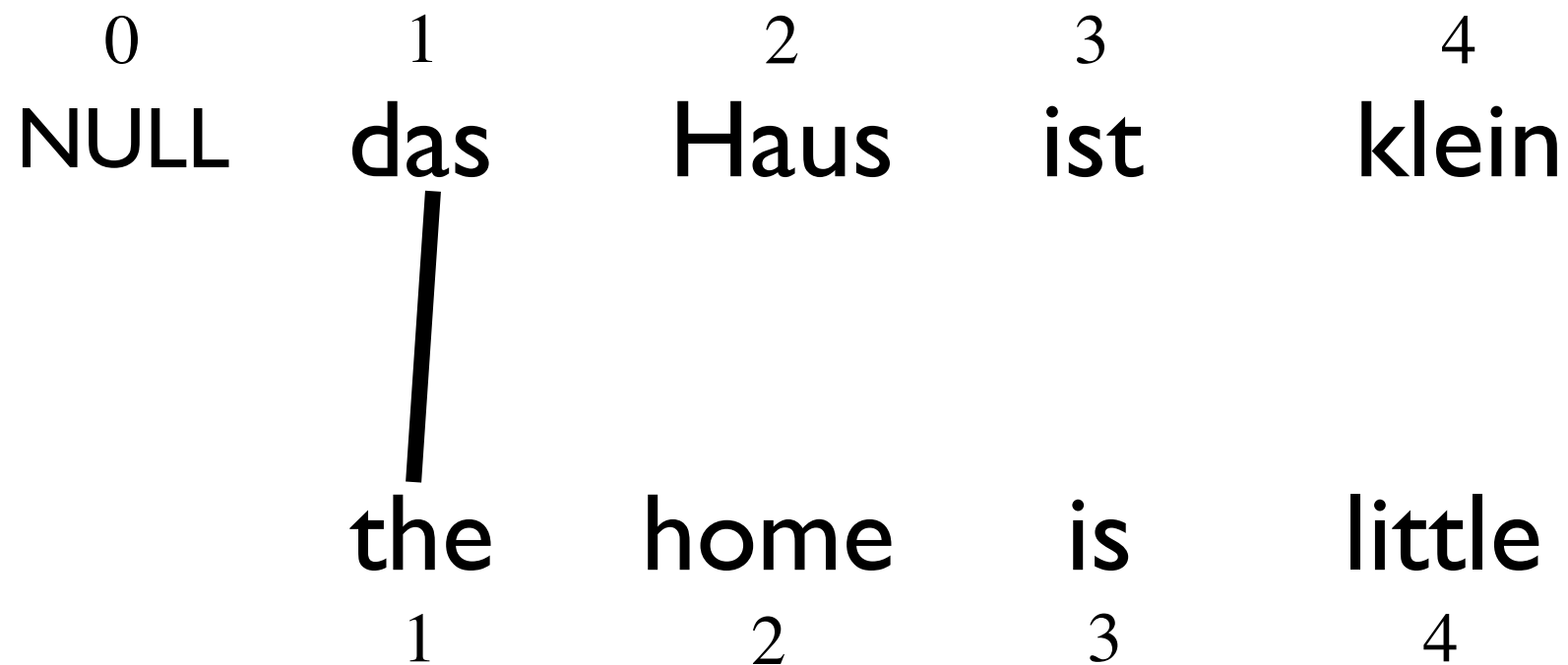
Finding the Viterbi Alignment



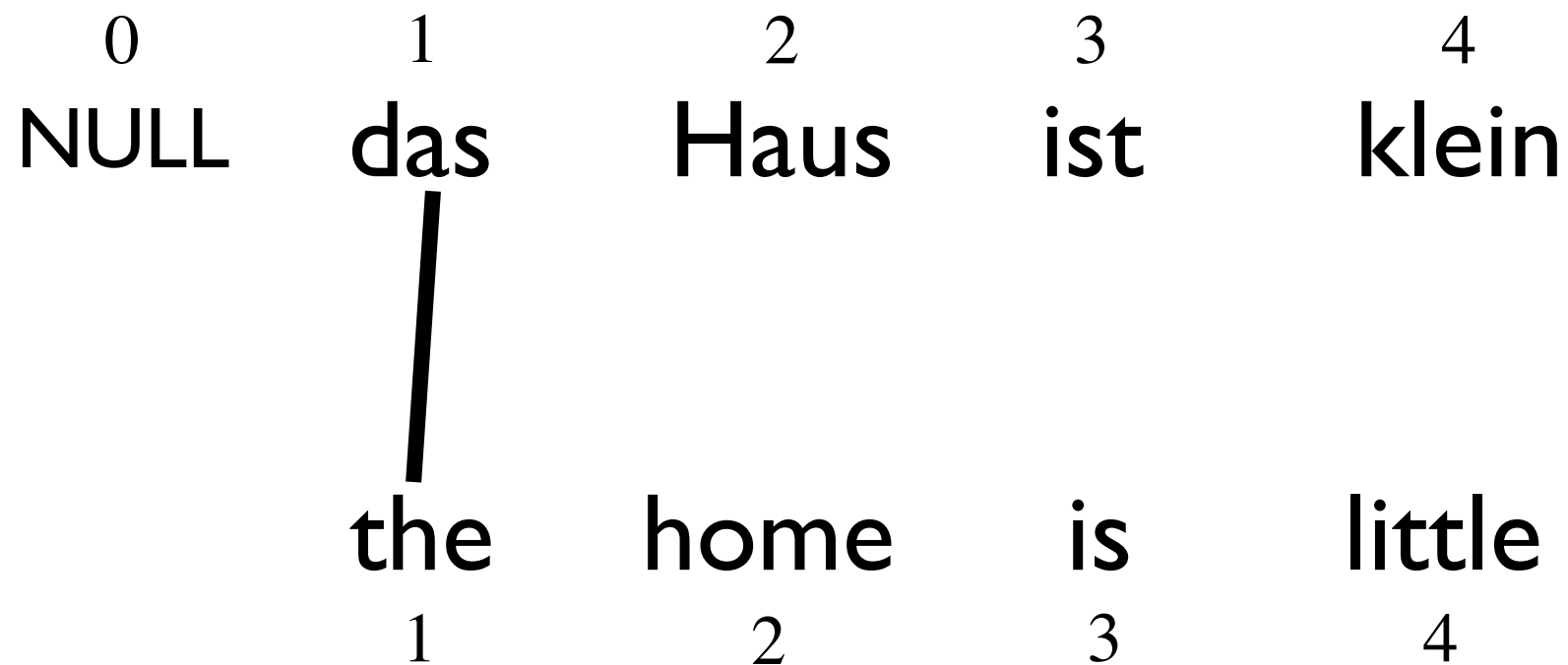
Finding the Viterbi Alignment



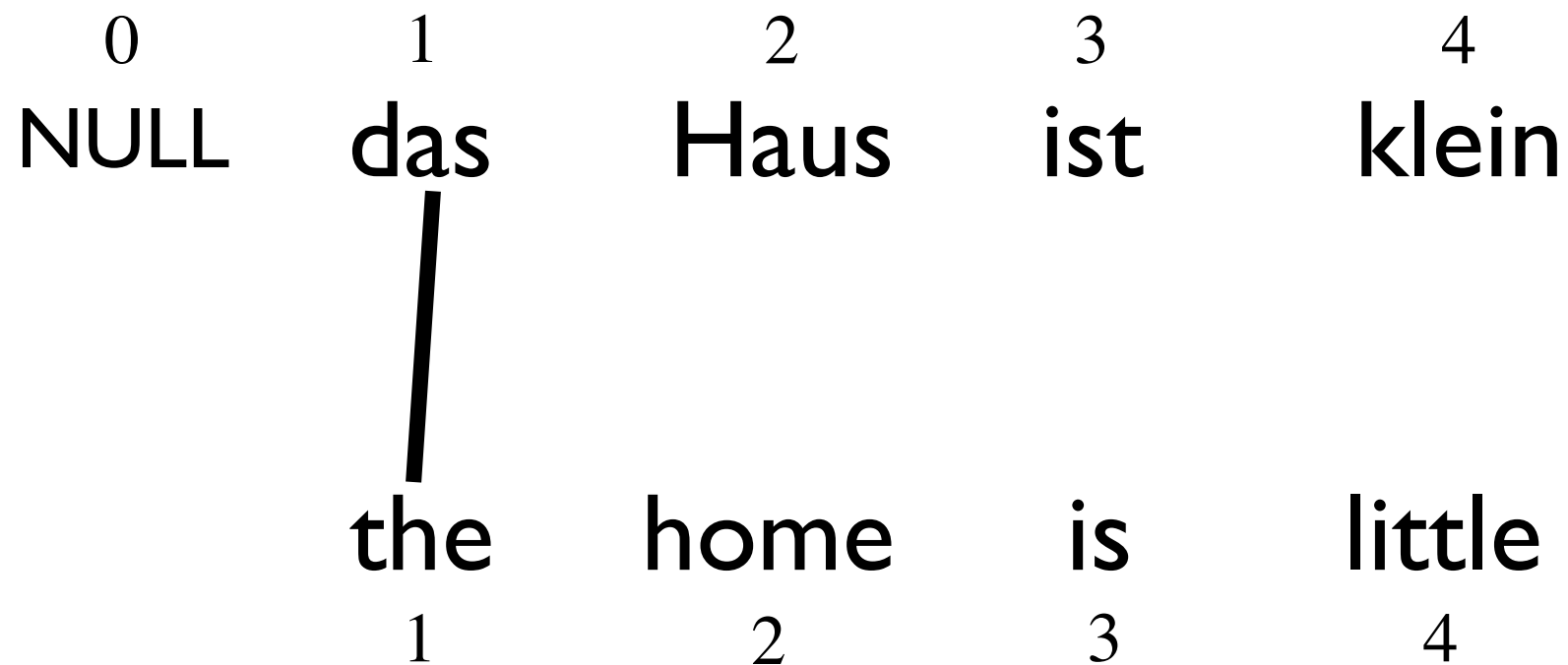
Finding the Viterbi Alignment



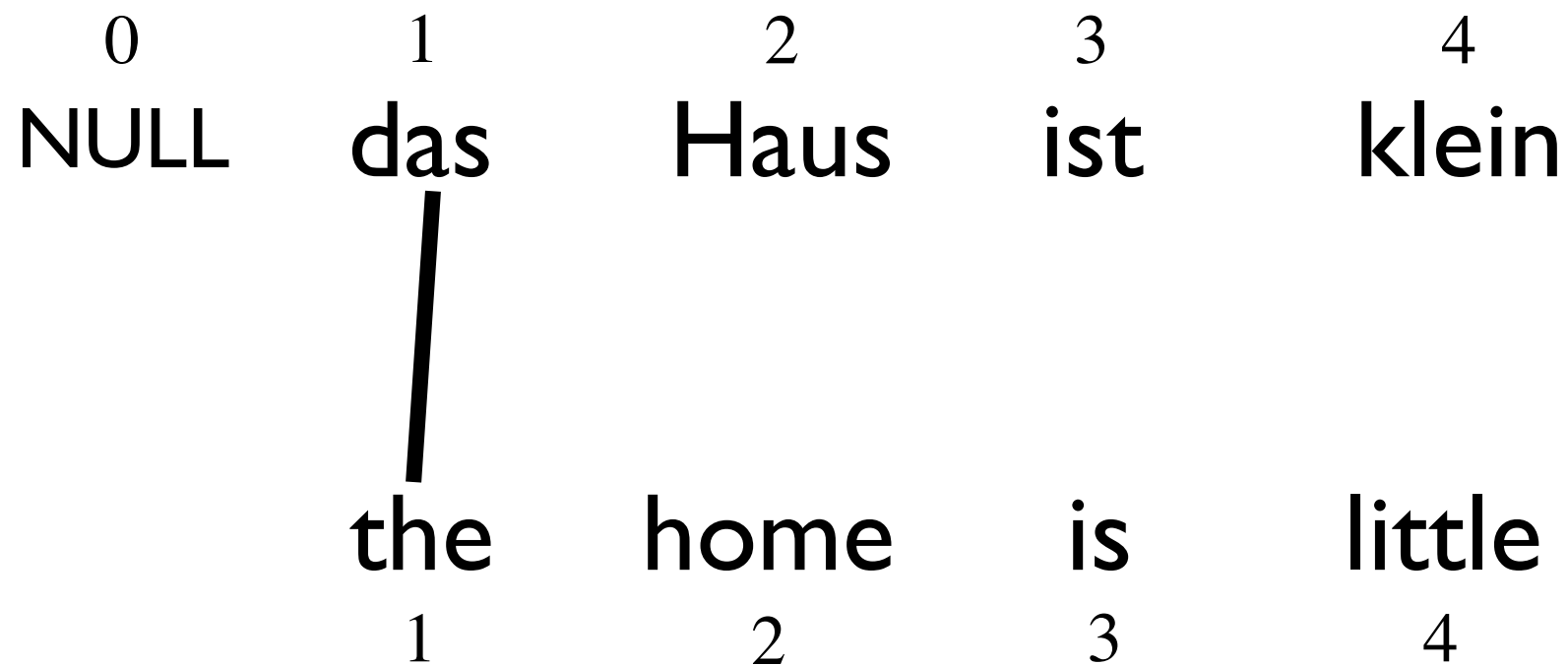
Finding the Viterbi Alignment



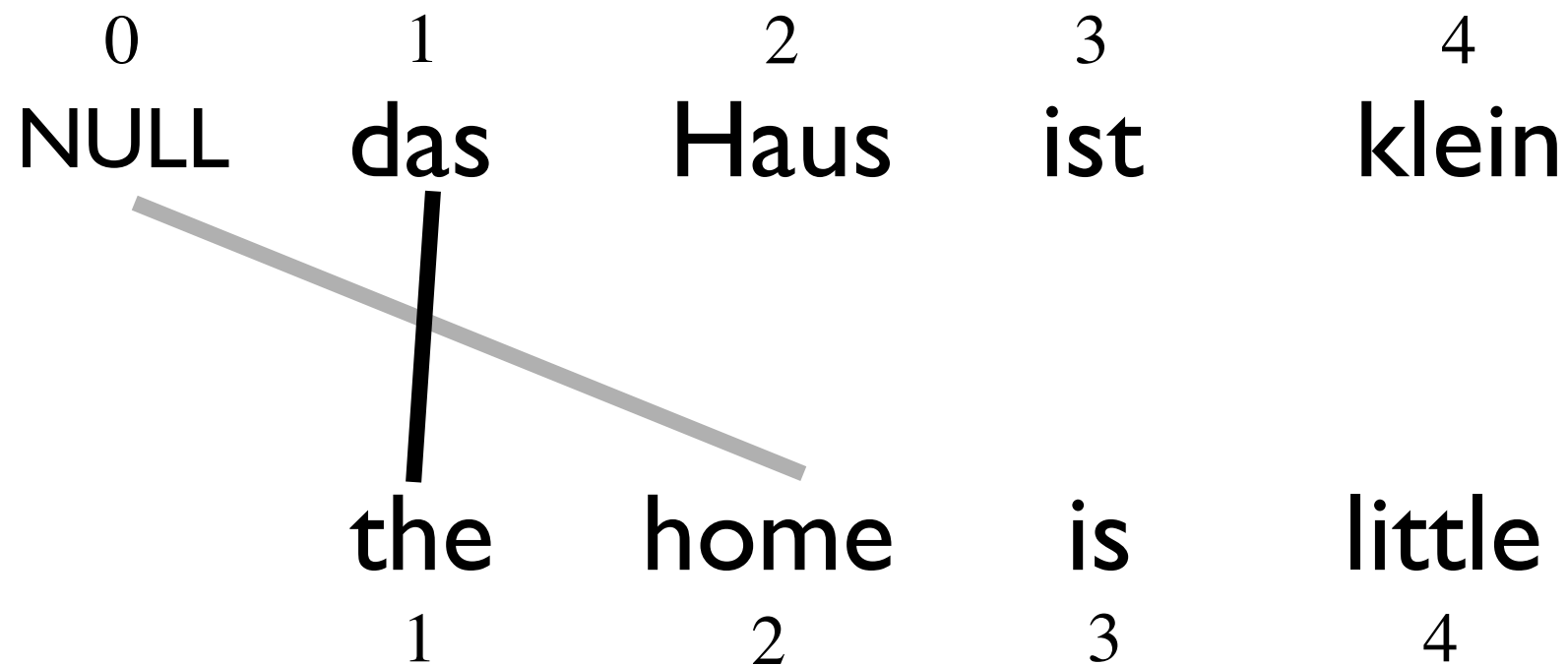
Finding the Viterbi Alignment



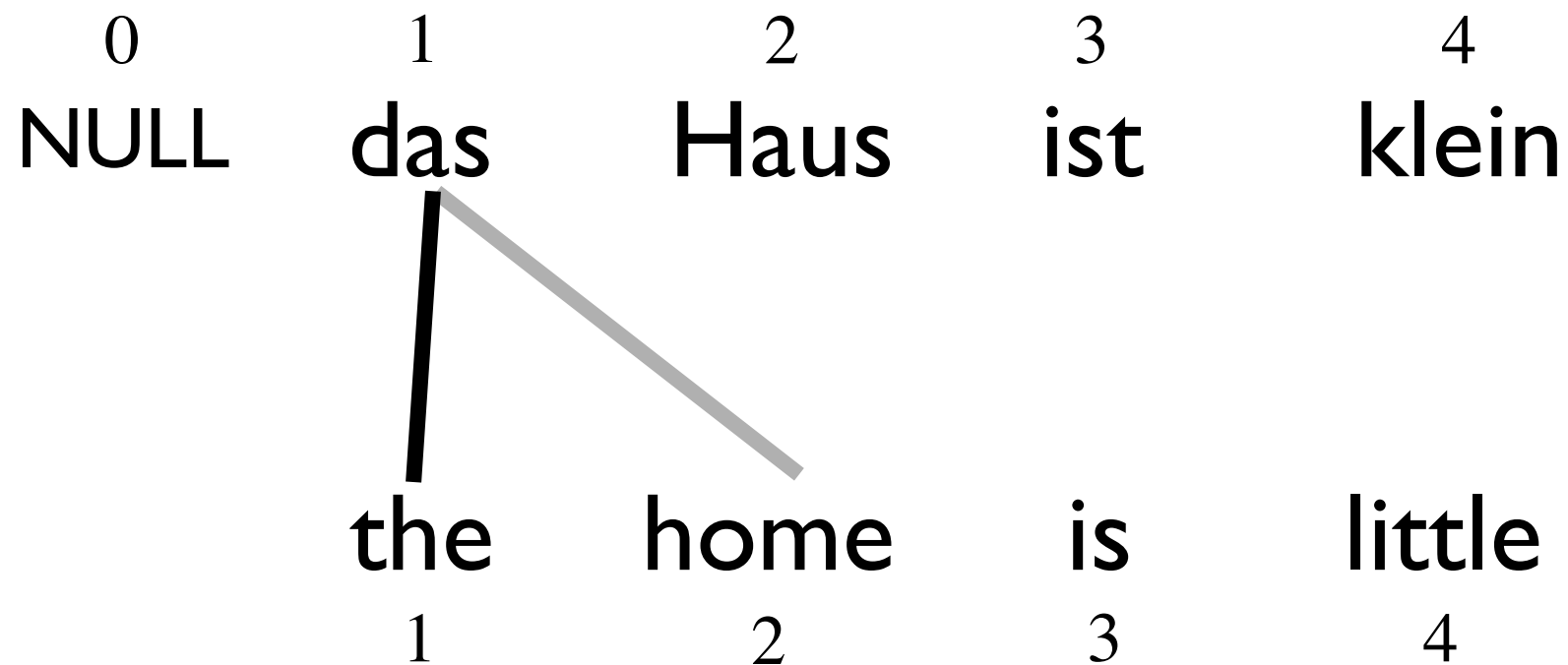
Finding the Viterbi Alignment



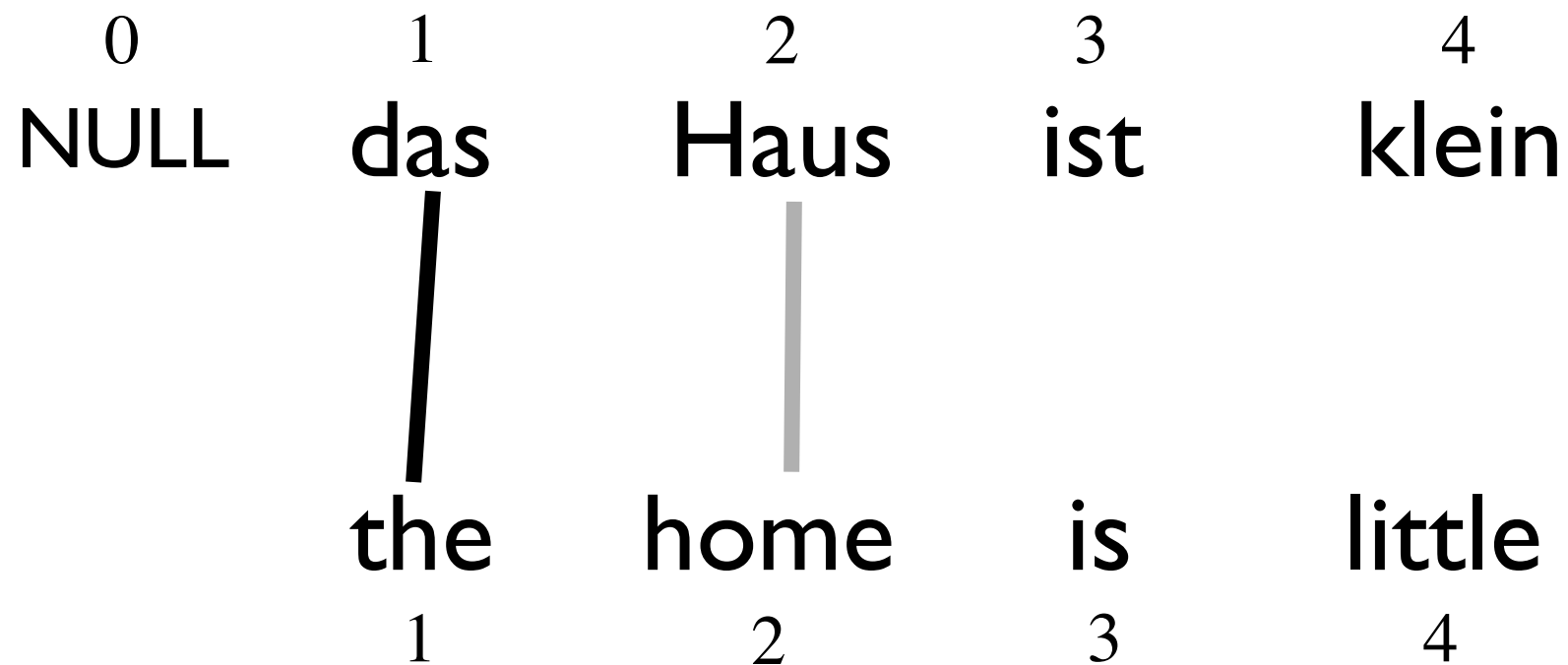
Finding the Viterbi Alignment



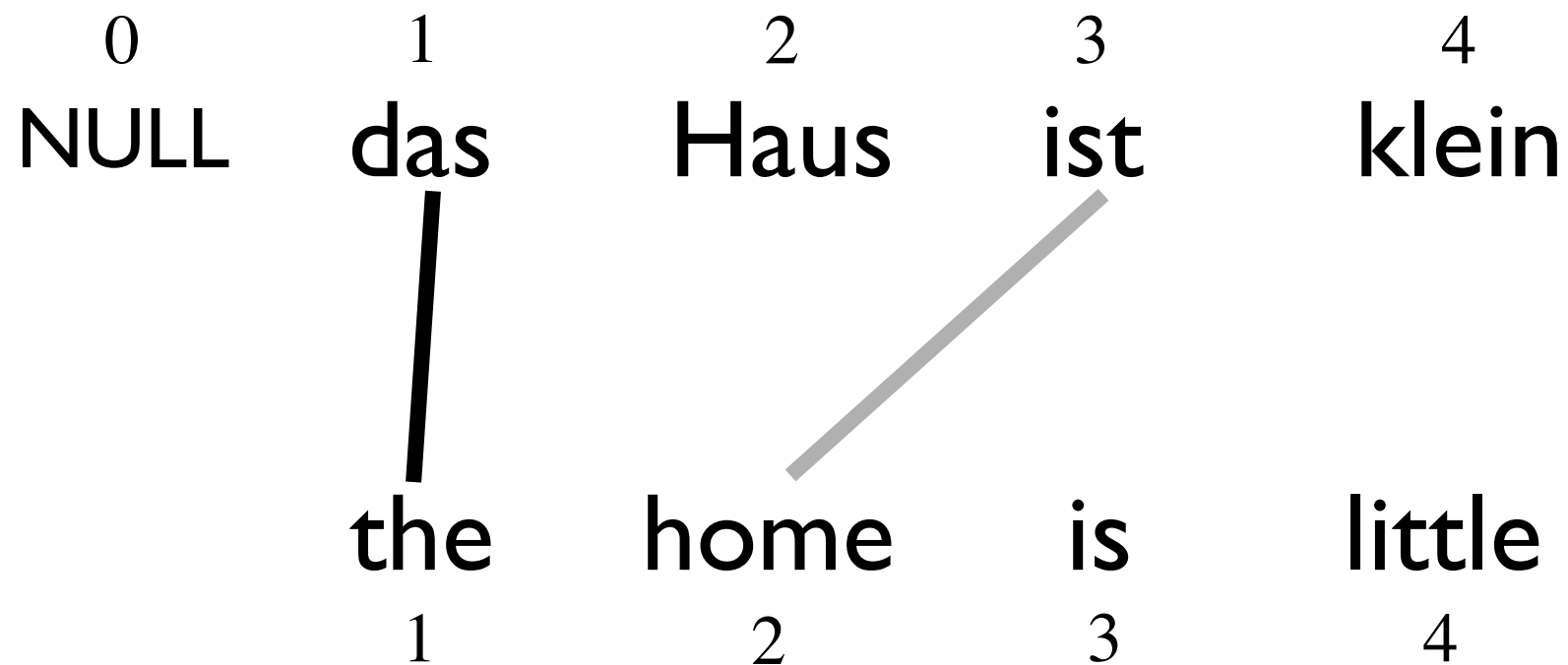
Finding the Viterbi Alignment



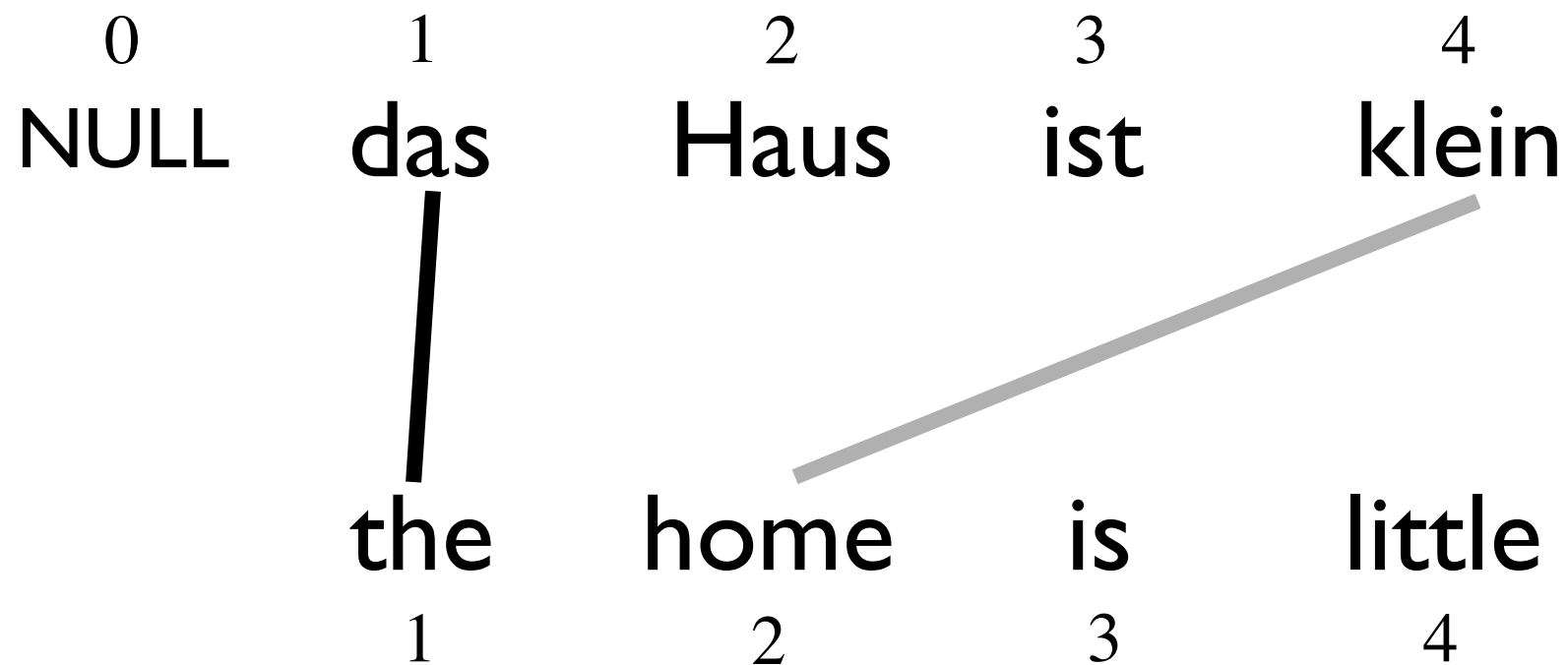
Finding the Viterbi Alignment



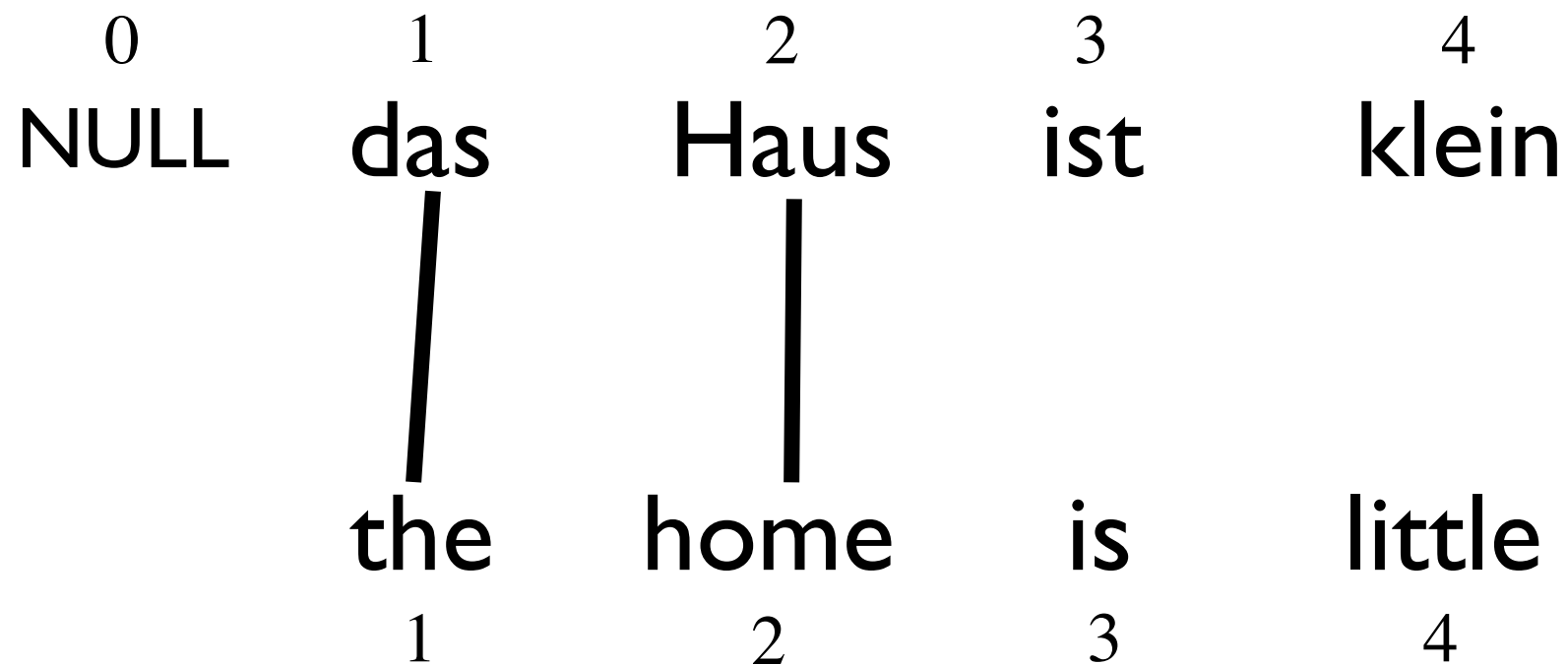
Finding the Viterbi Alignment



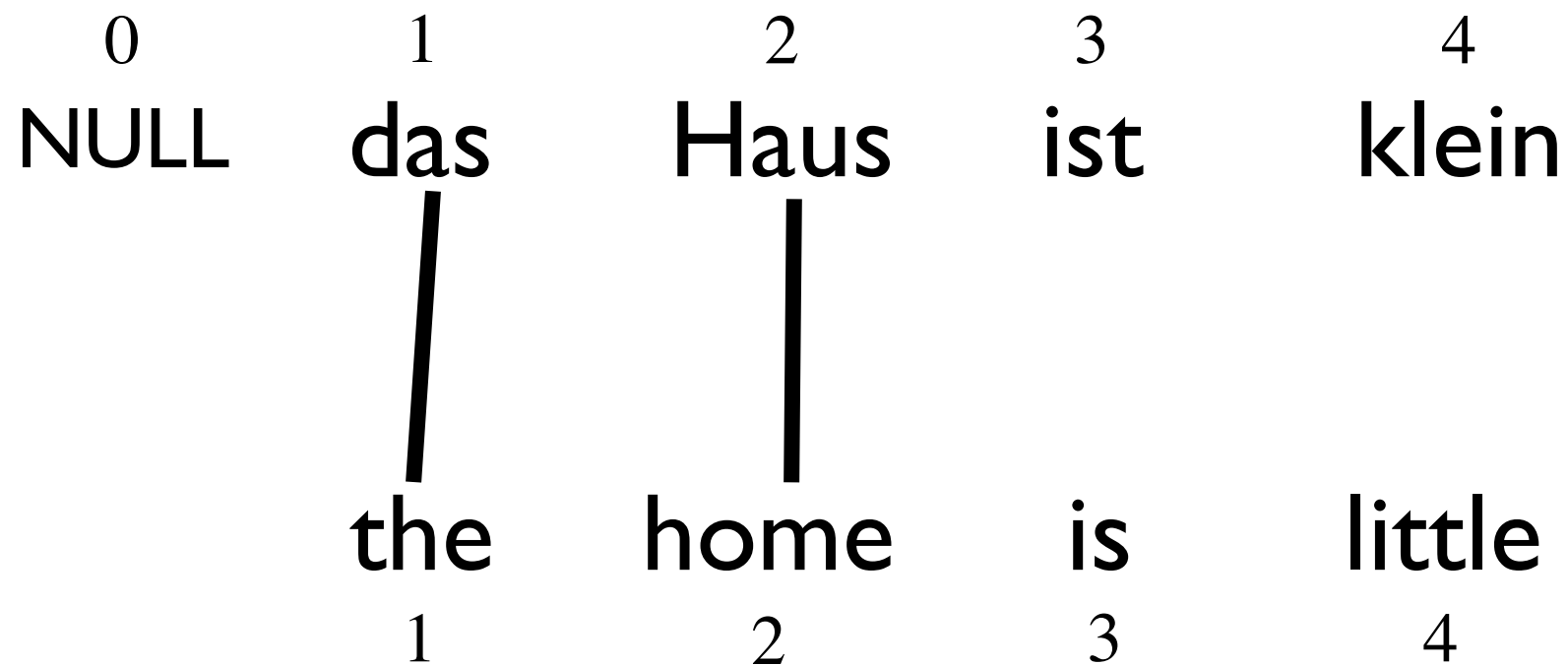
Finding the Viterbi Alignment



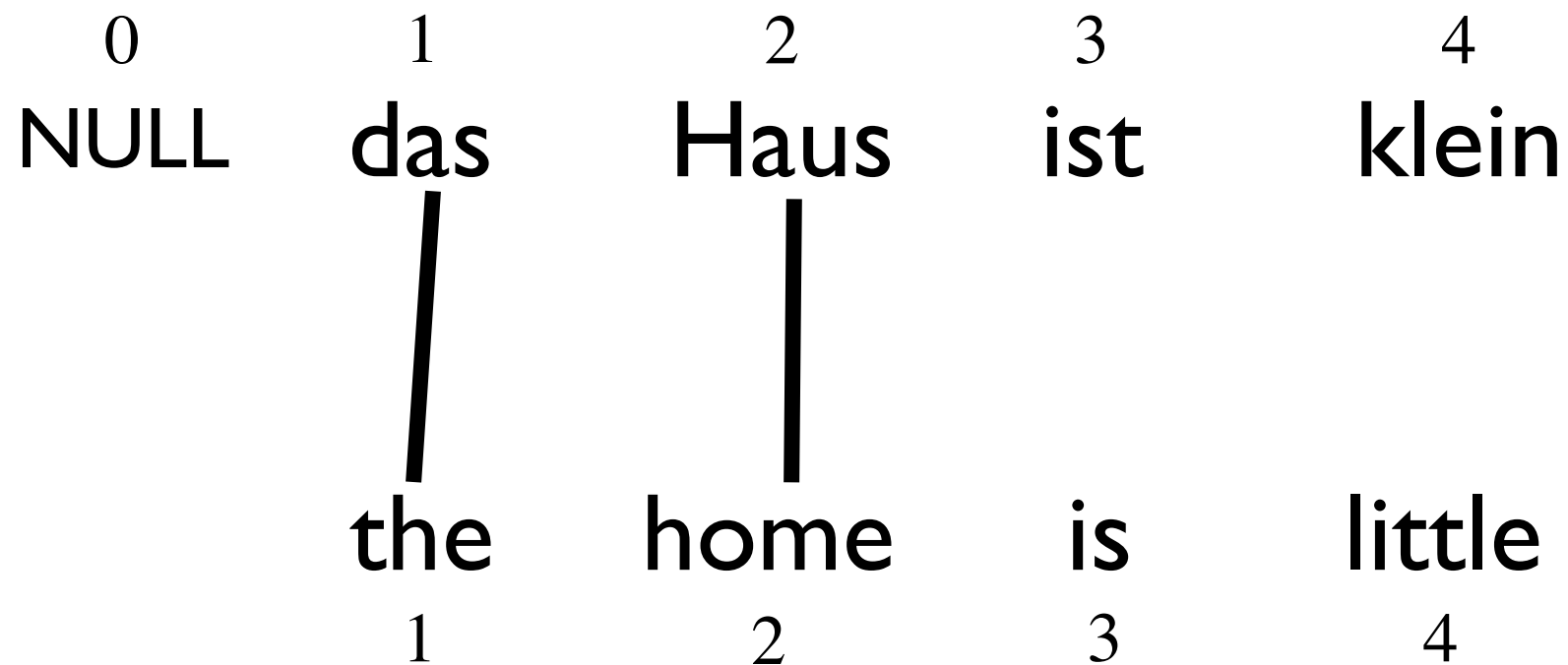
Finding the Viterbi Alignment



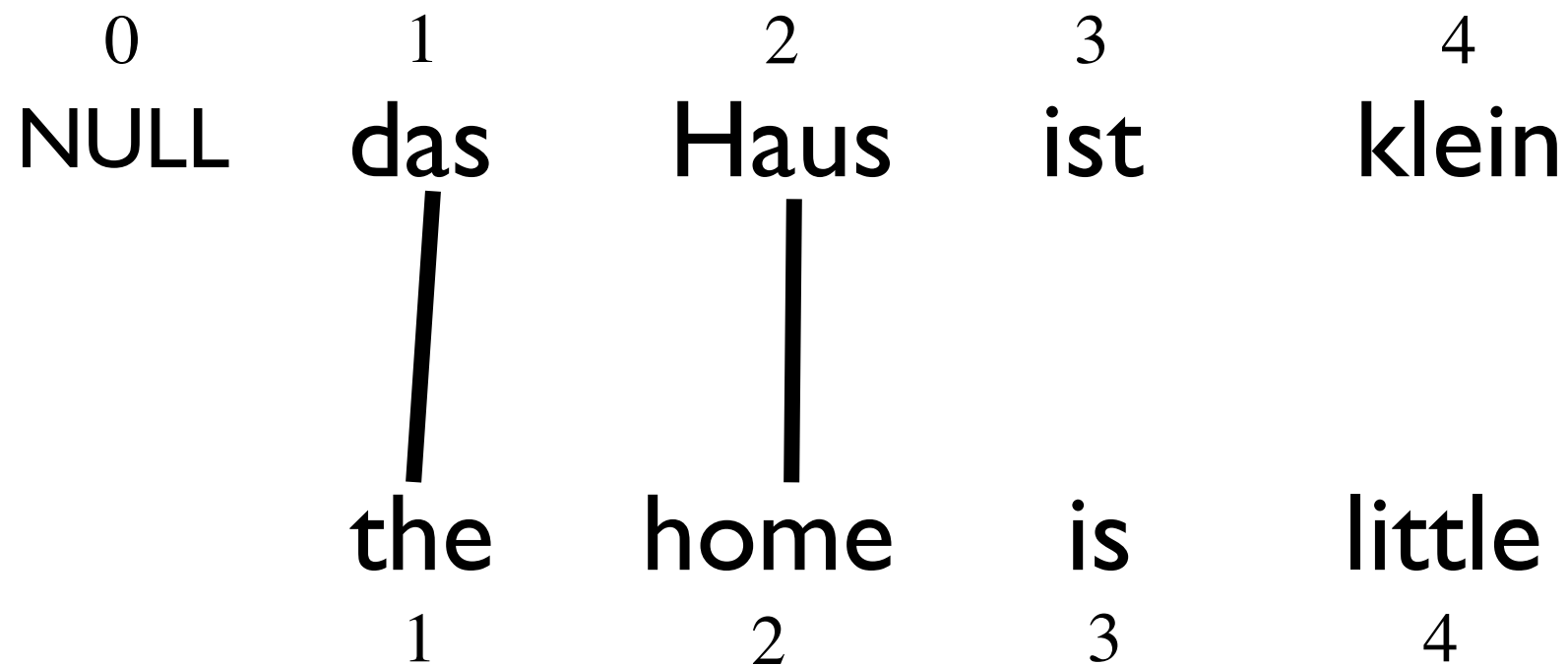
Finding the Viterbi Alignment



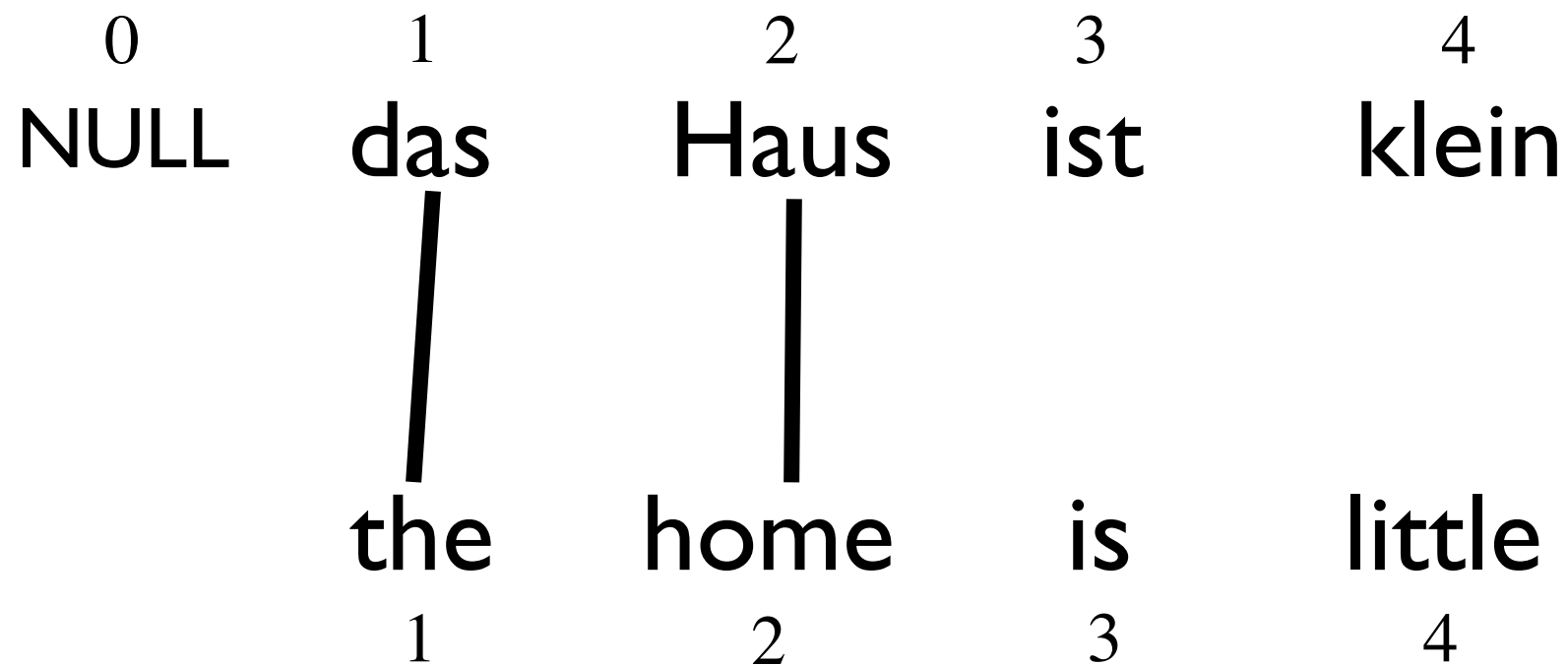
Finding the Viterbi Alignment



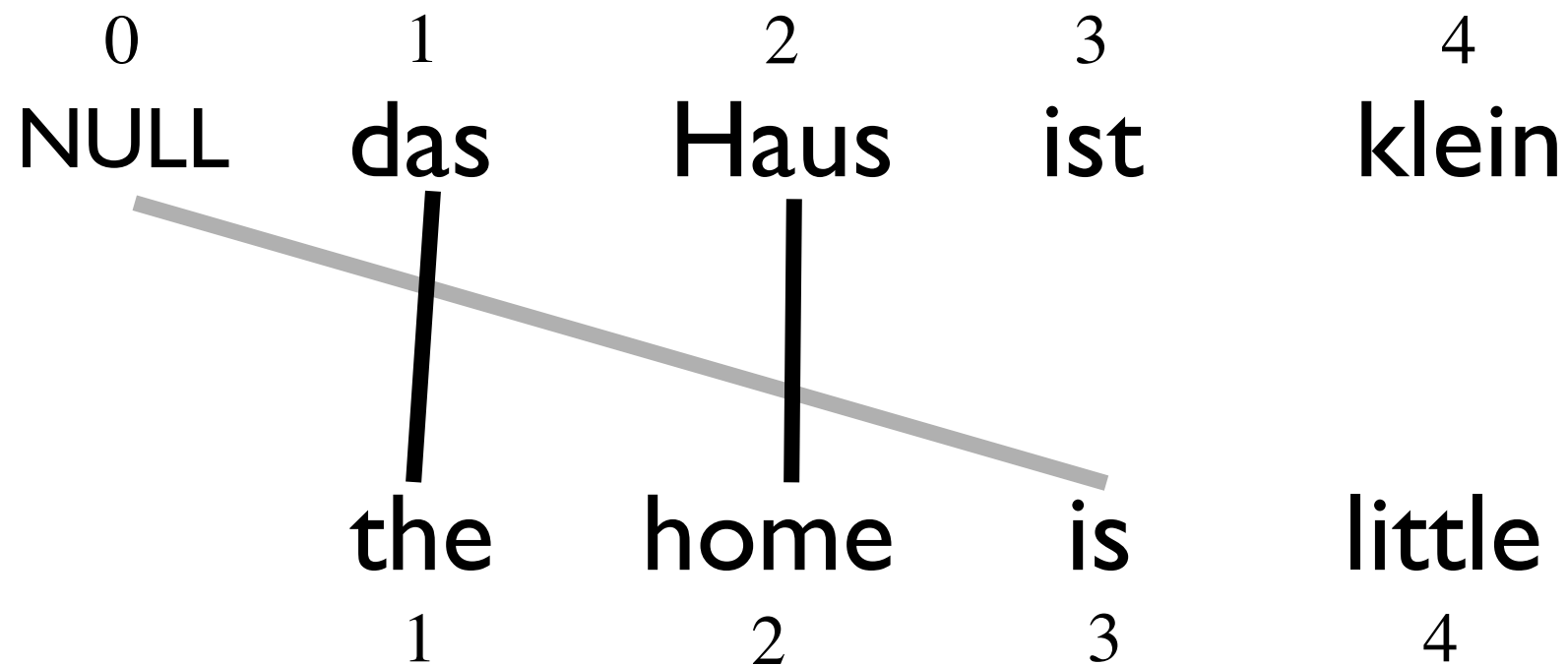
Finding the Viterbi Alignment



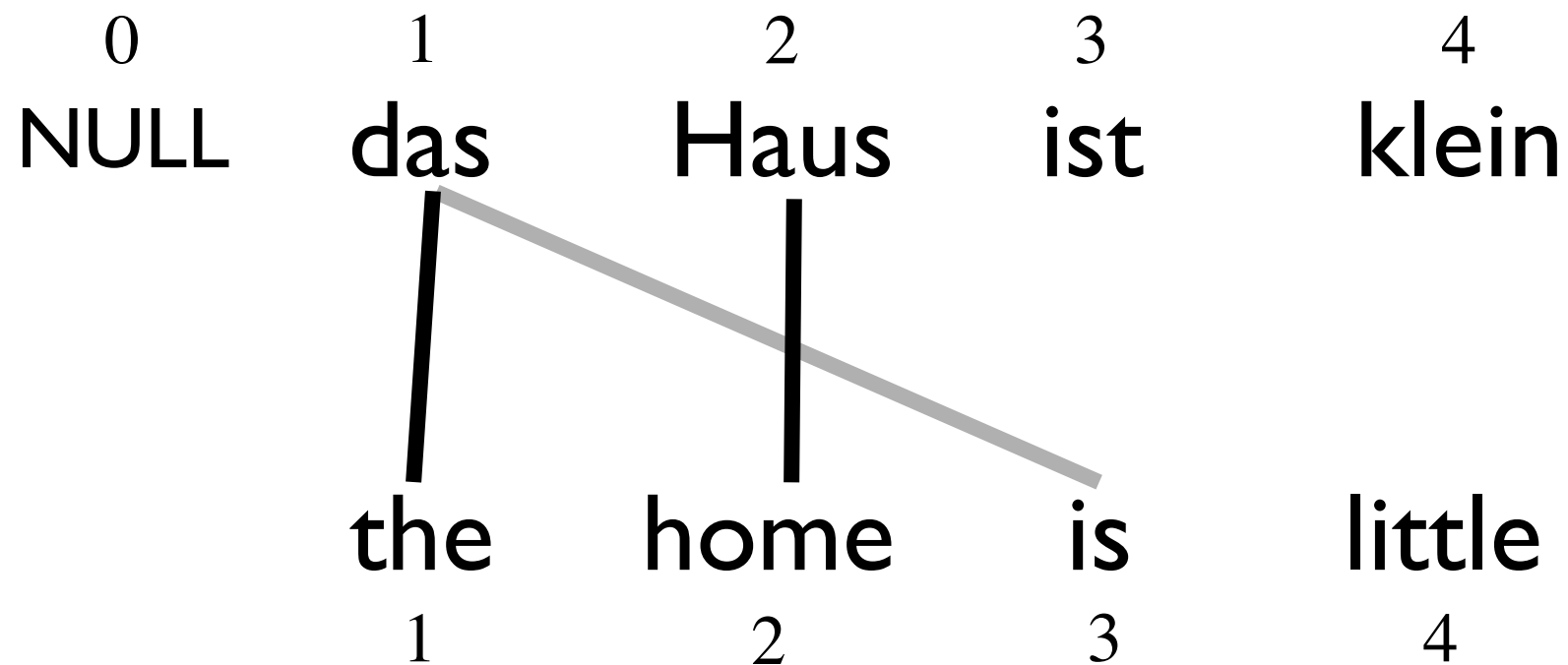
Finding the Viterbi Alignment



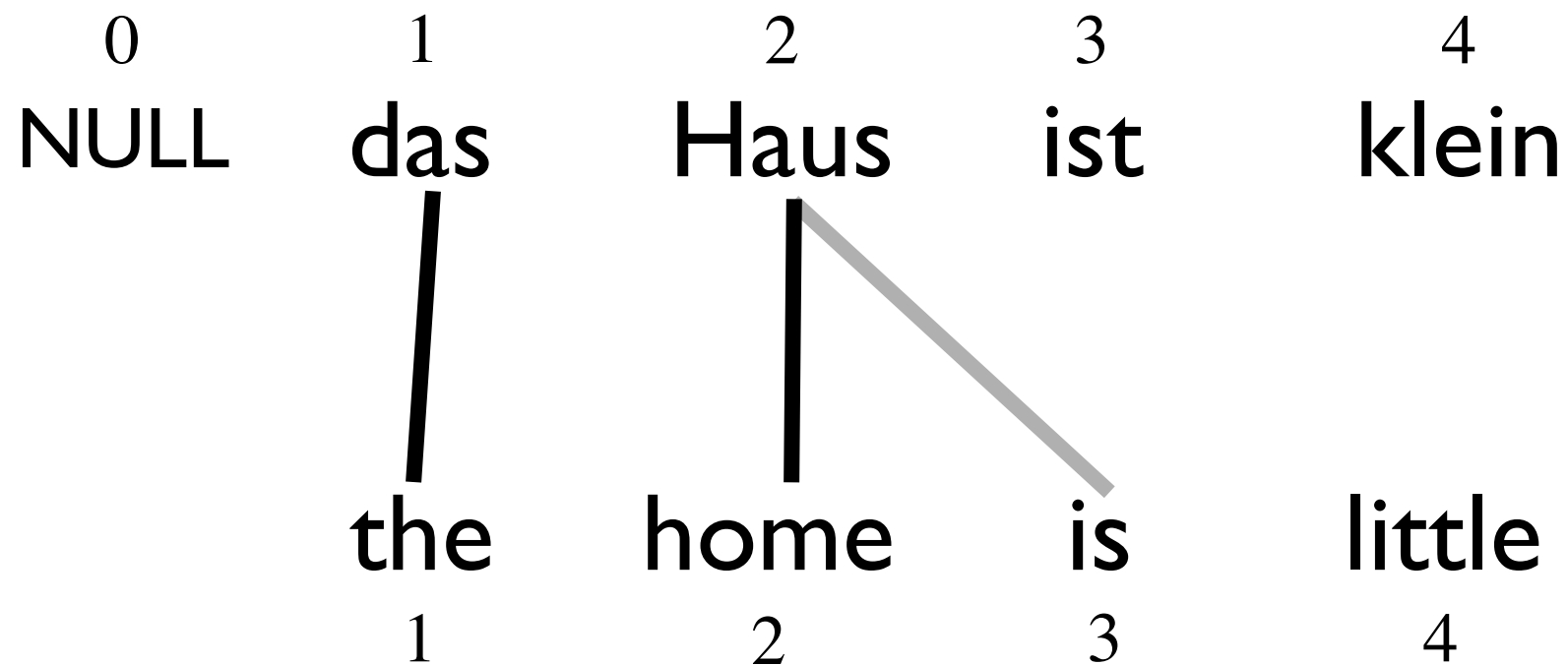
Finding the Viterbi Alignment



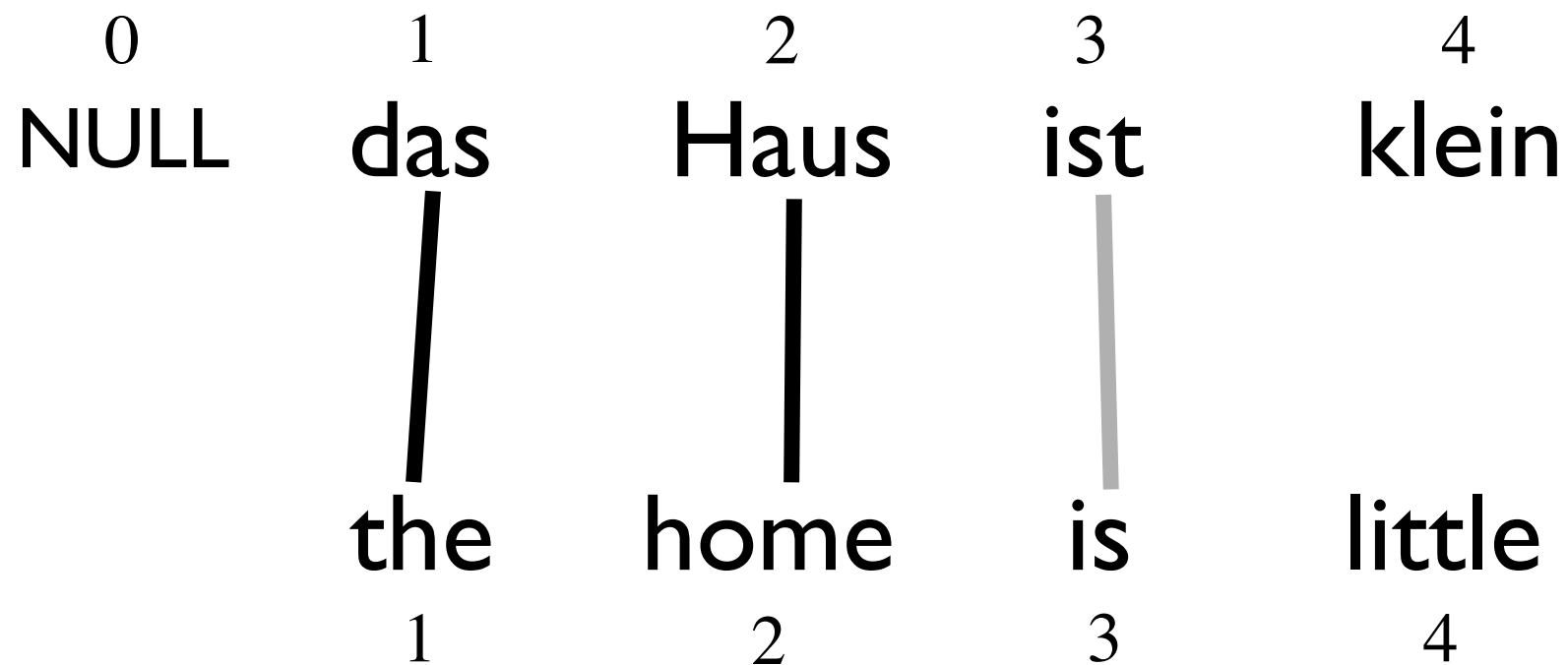
Finding the Viterbi Alignment



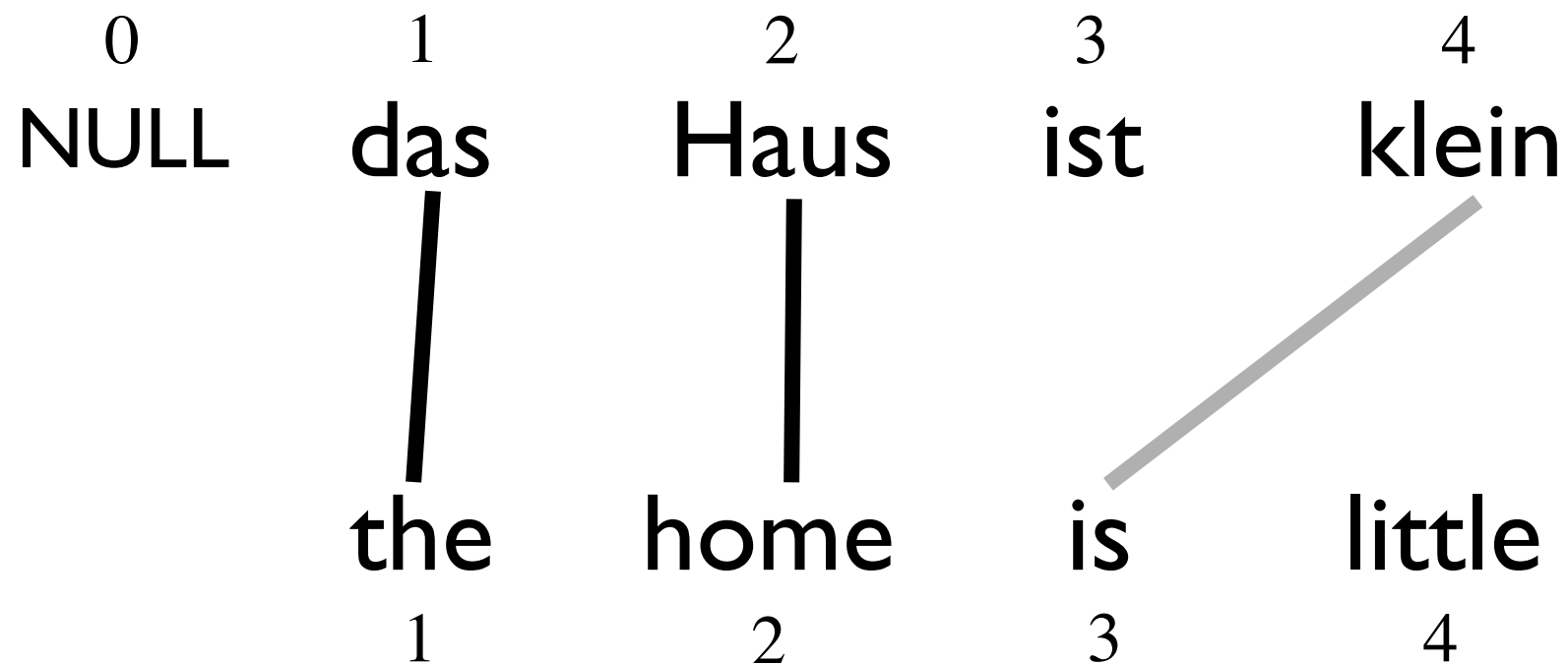
Finding the Viterbi Alignment



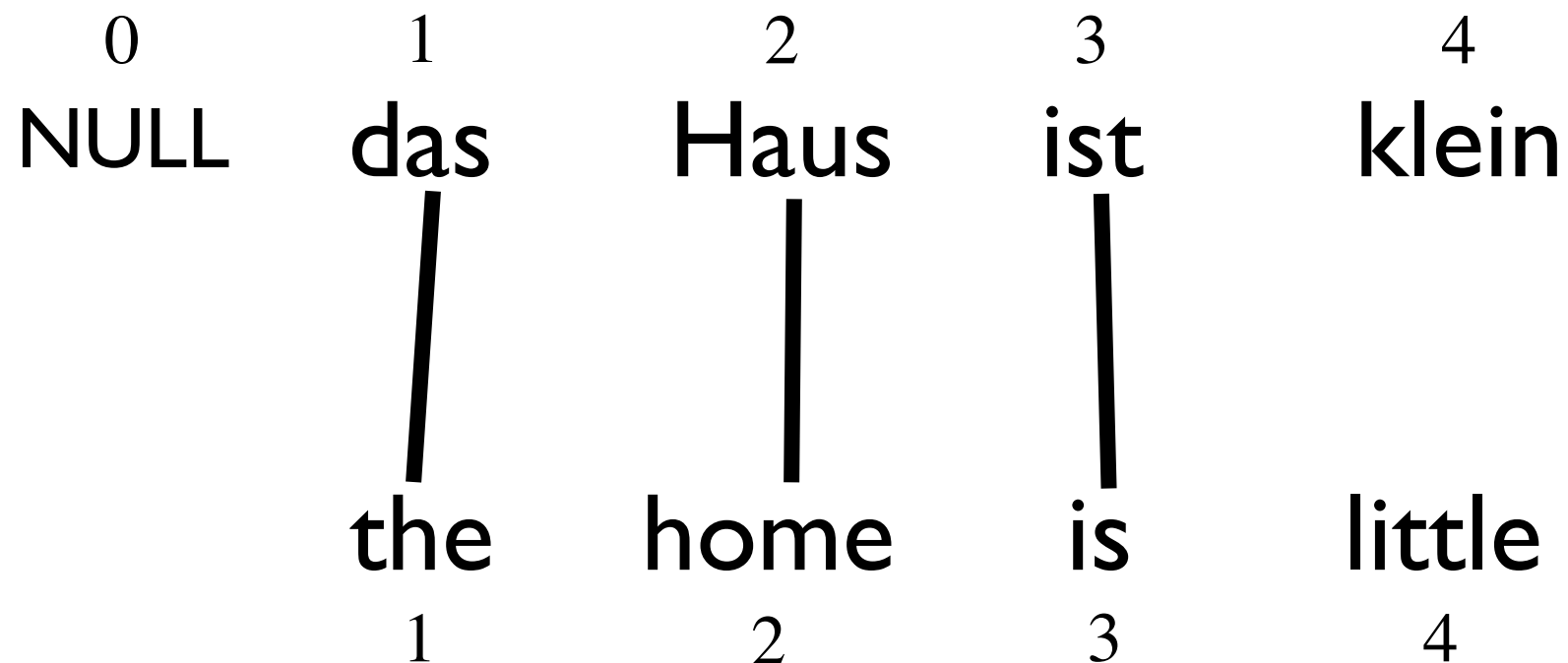
Finding the Viterbi Alignment



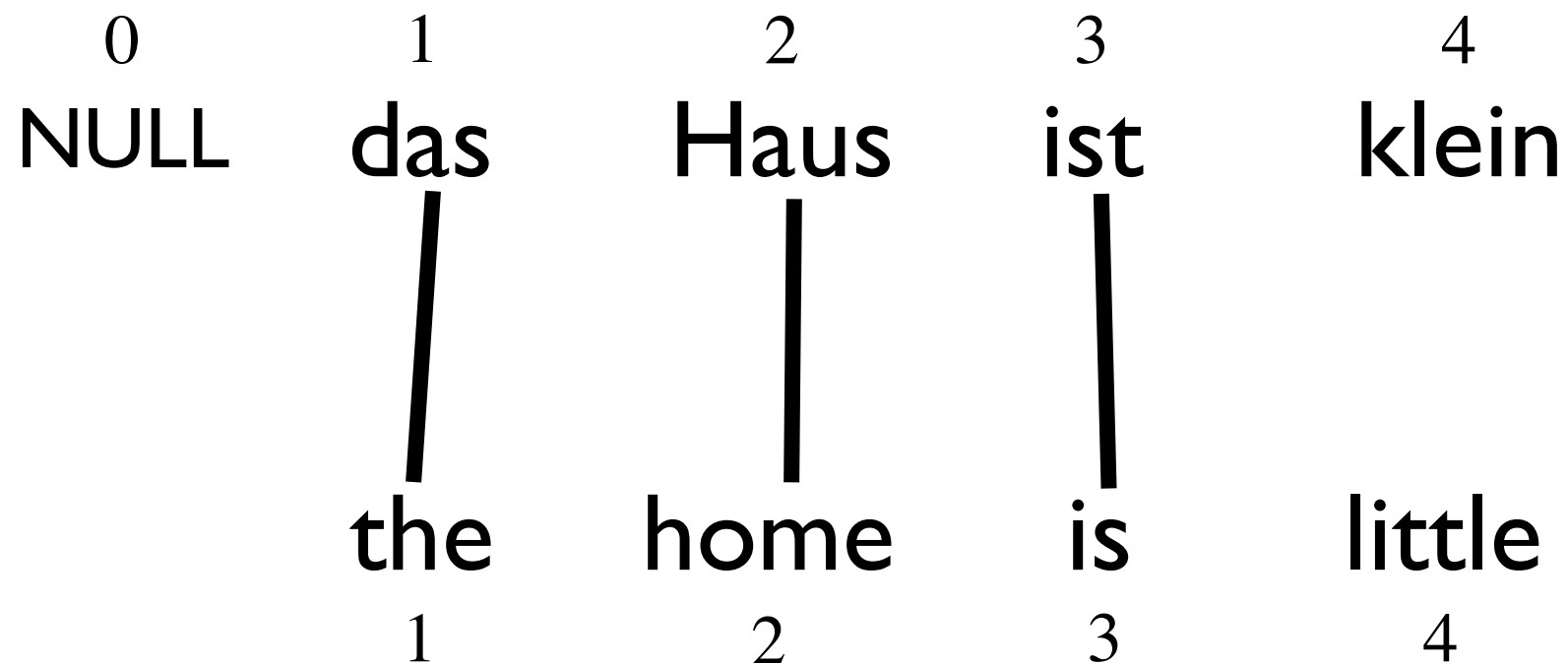
Finding the Viterbi Alignment



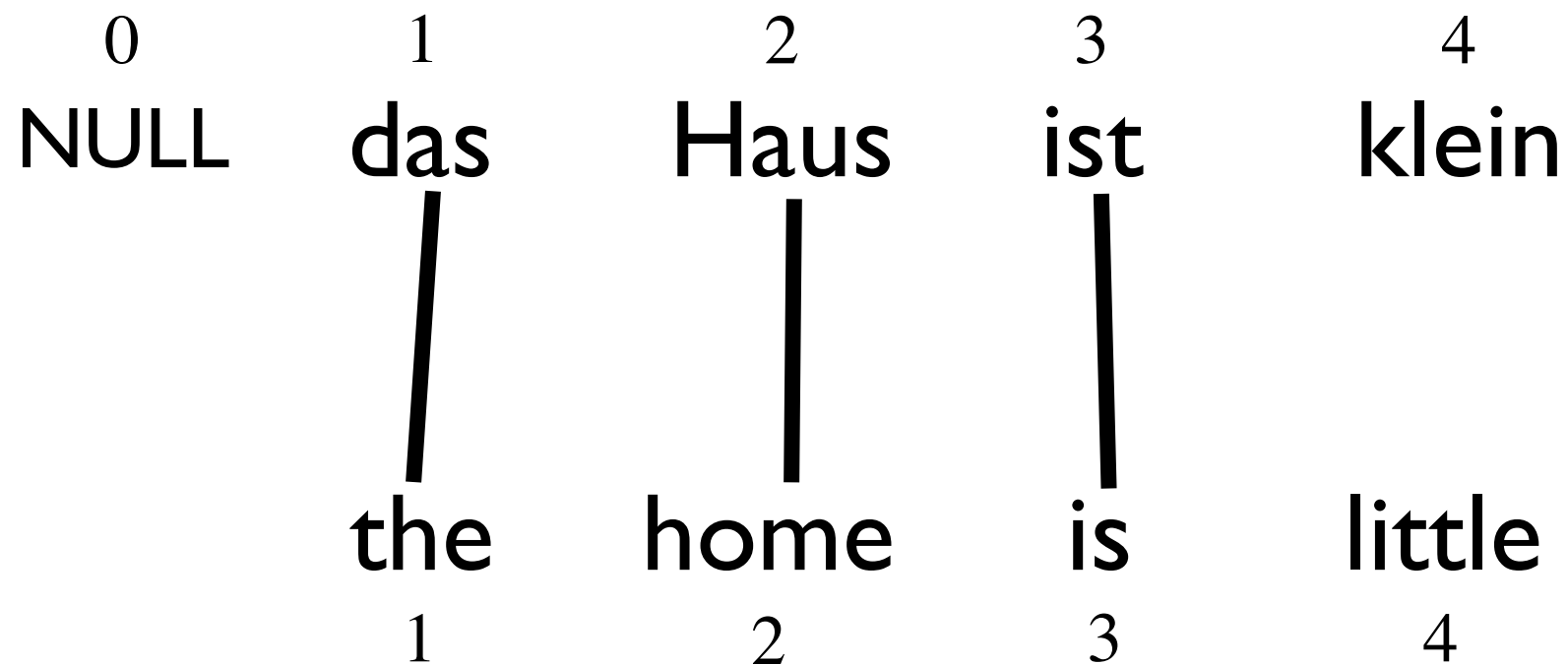
Finding the Viterbi Alignment



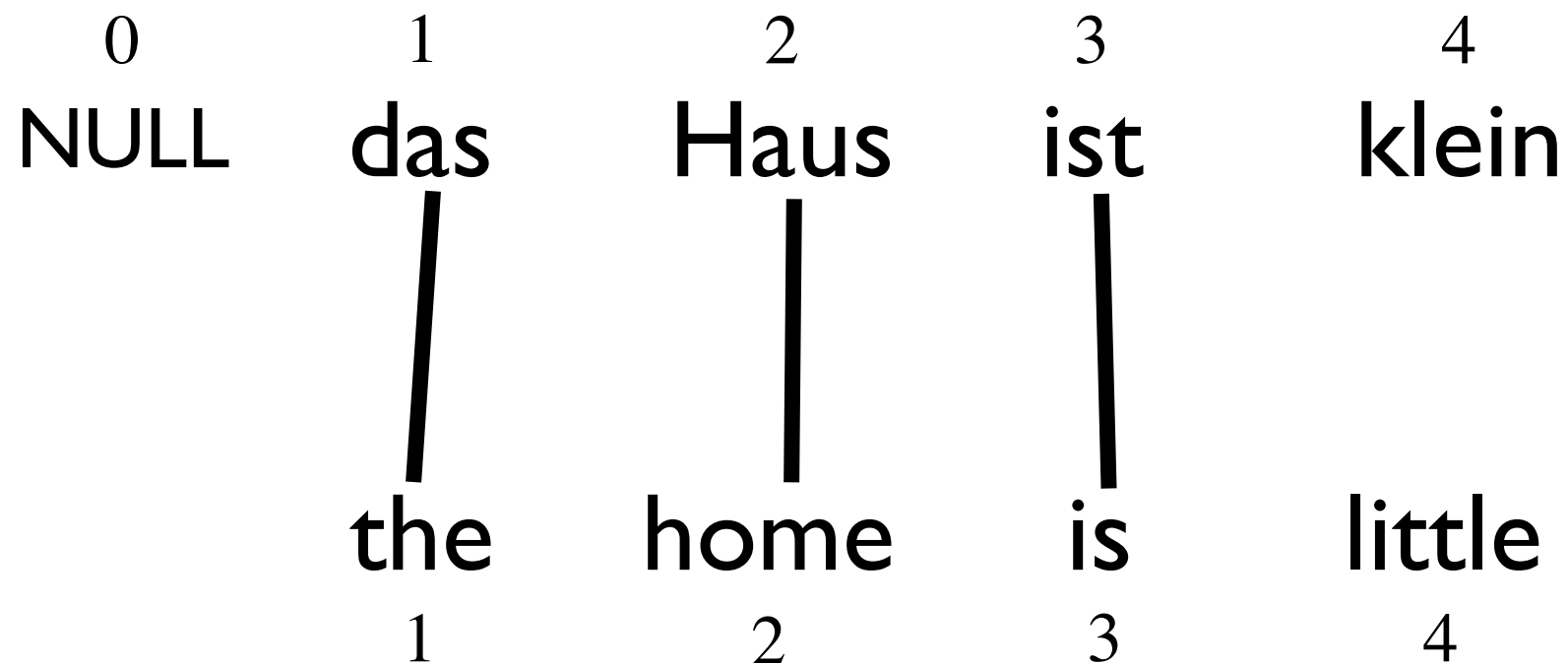
Finding the Viterbi Alignment



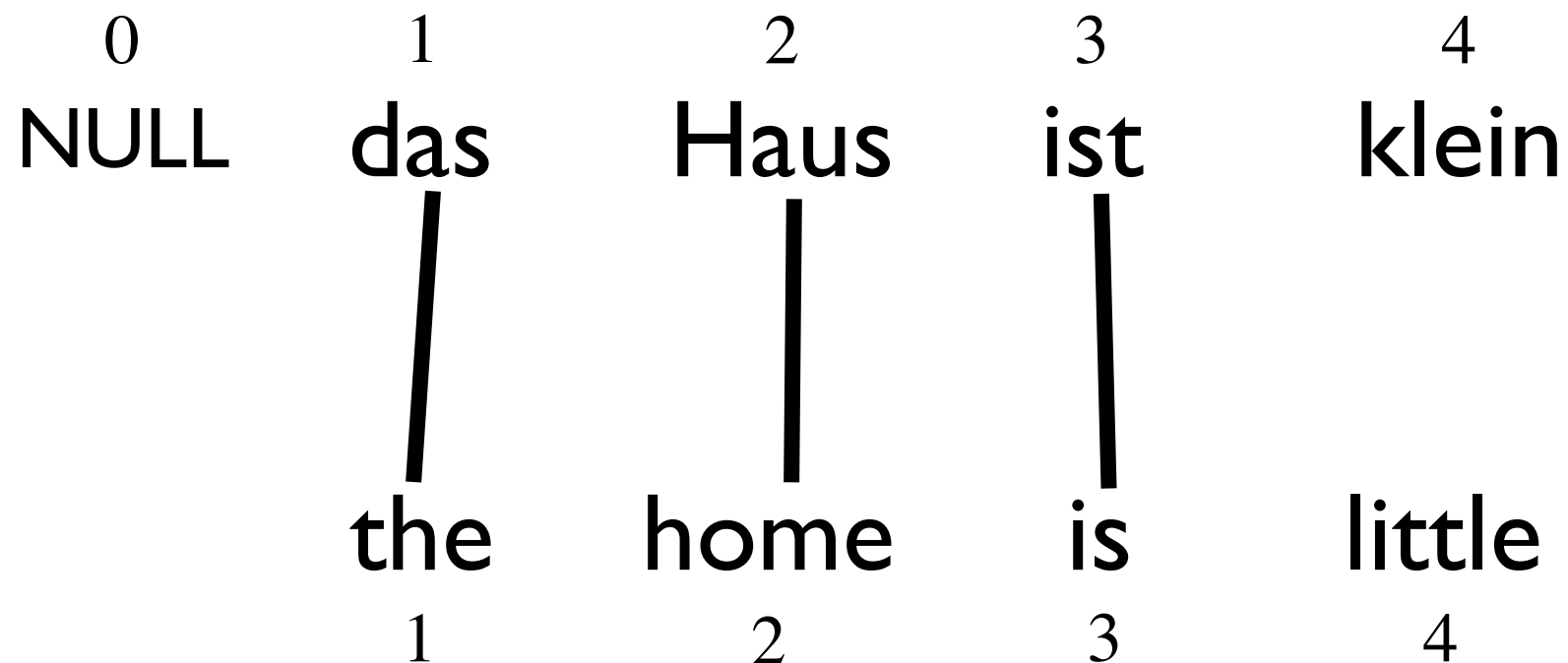
Finding the Viterbi Alignment



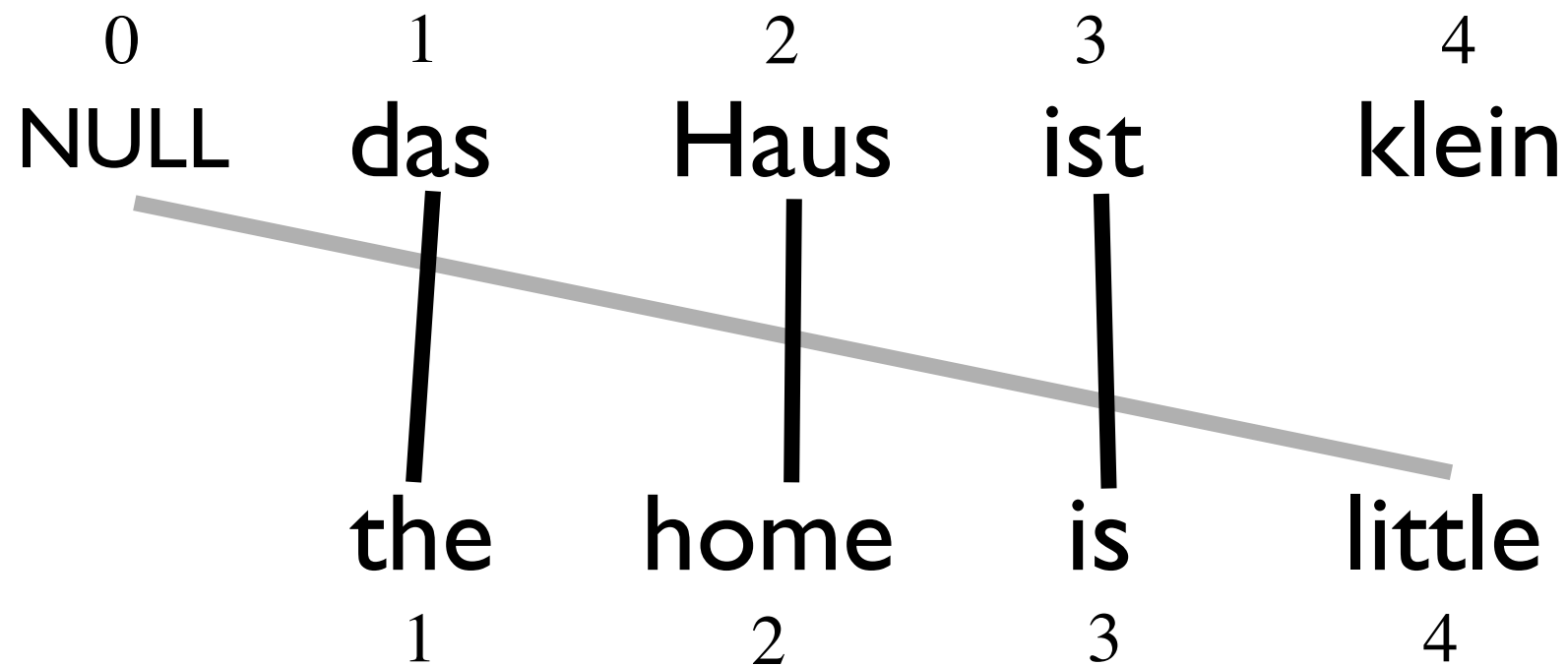
Finding the Viterbi Alignment



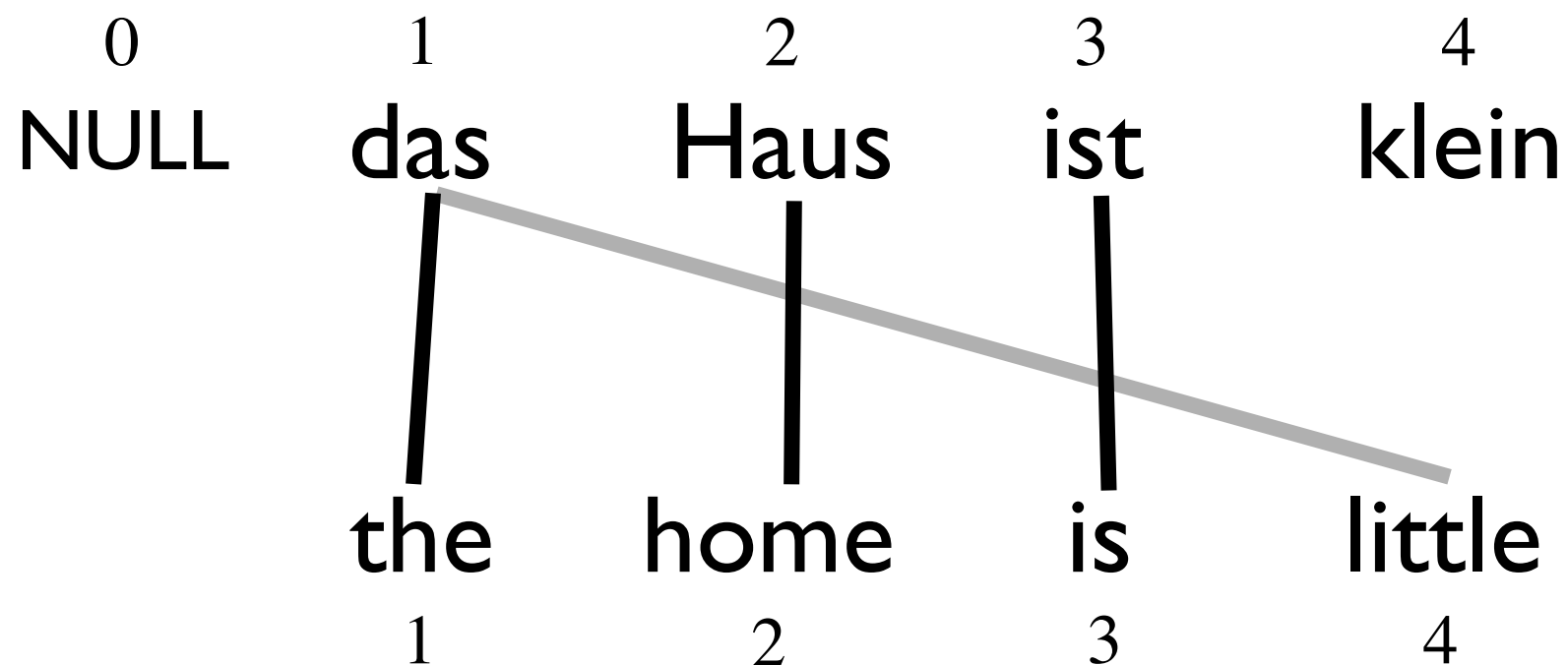
Finding the Viterbi Alignment



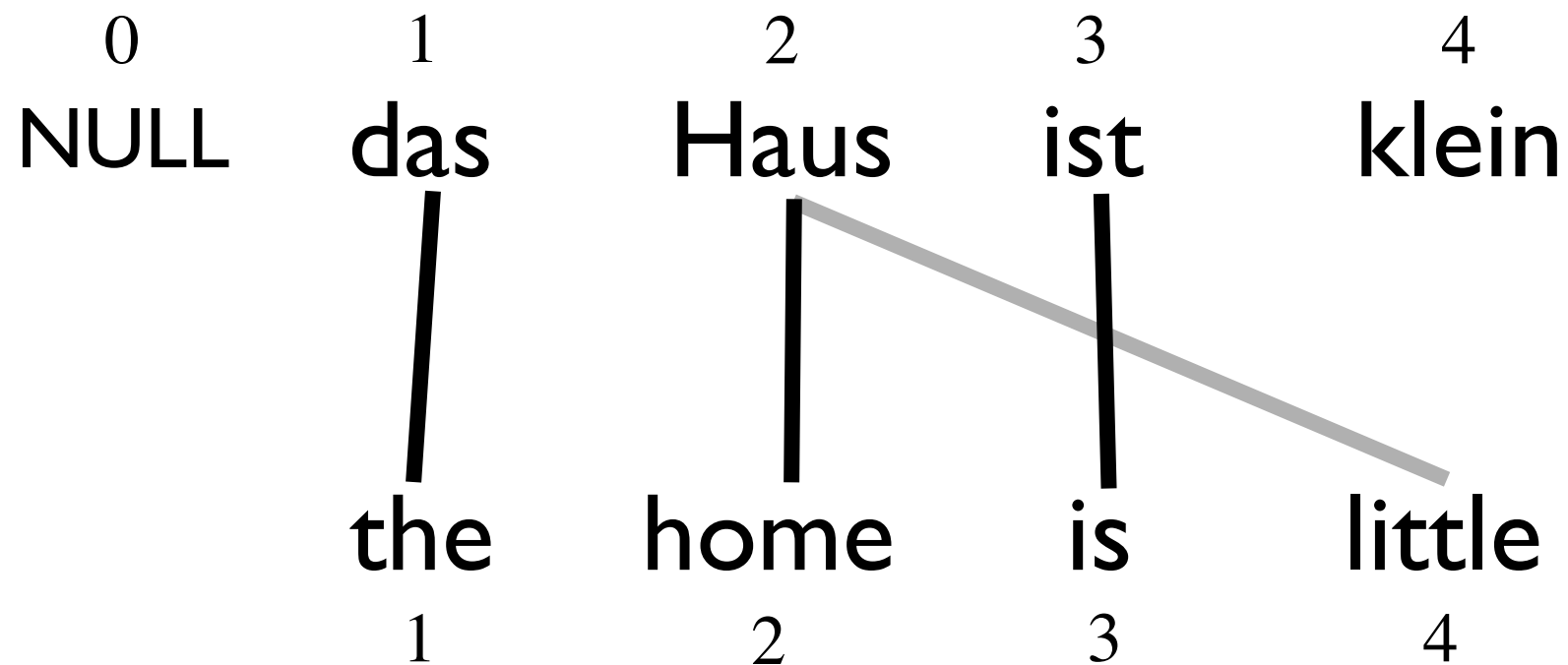
Finding the Viterbi Alignment



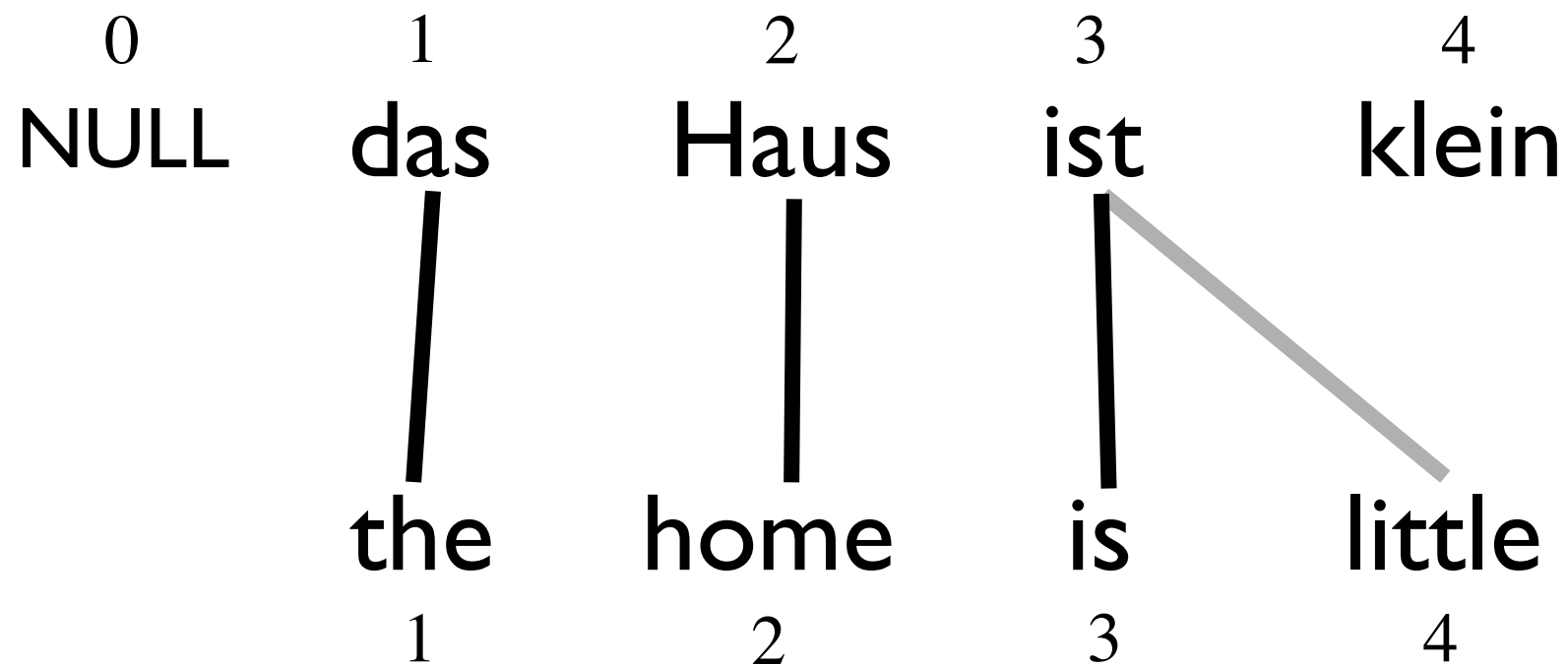
Finding the Viterbi Alignment



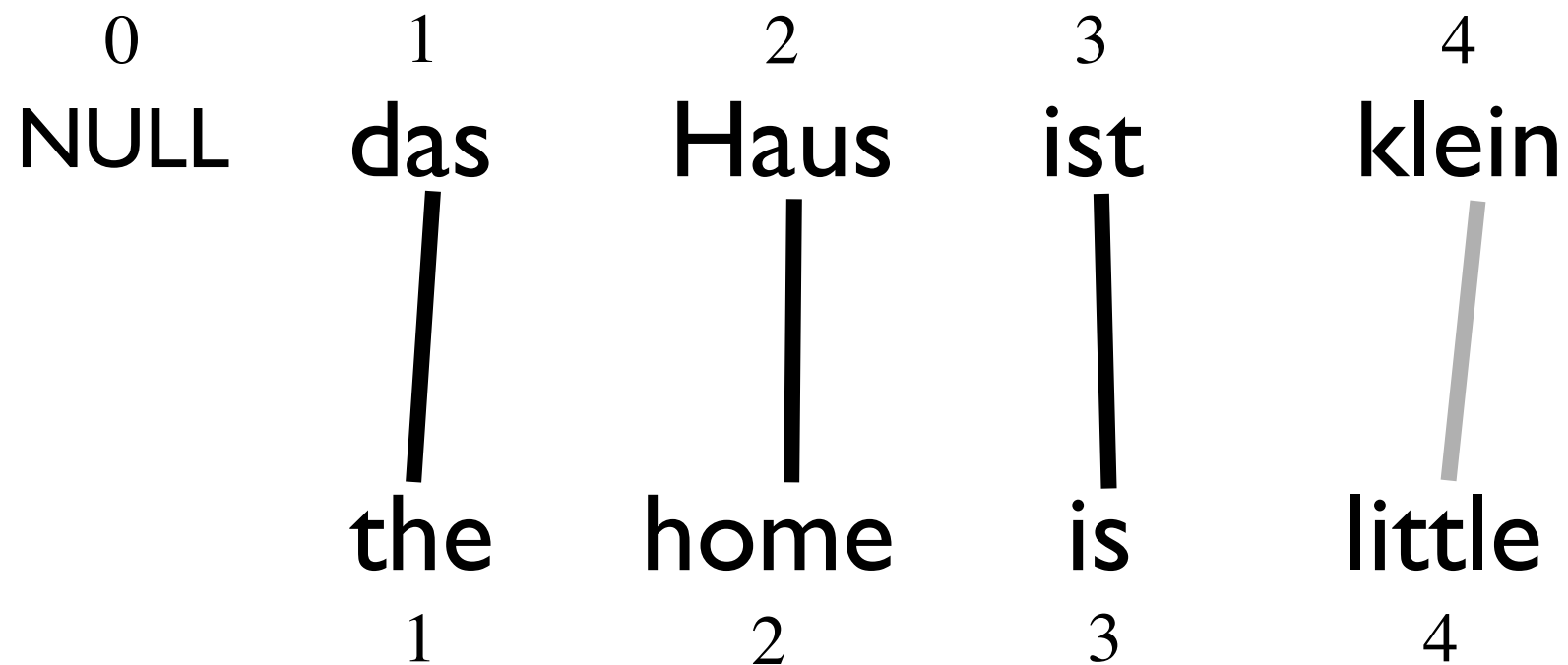
Finding the Viterbi Alignment



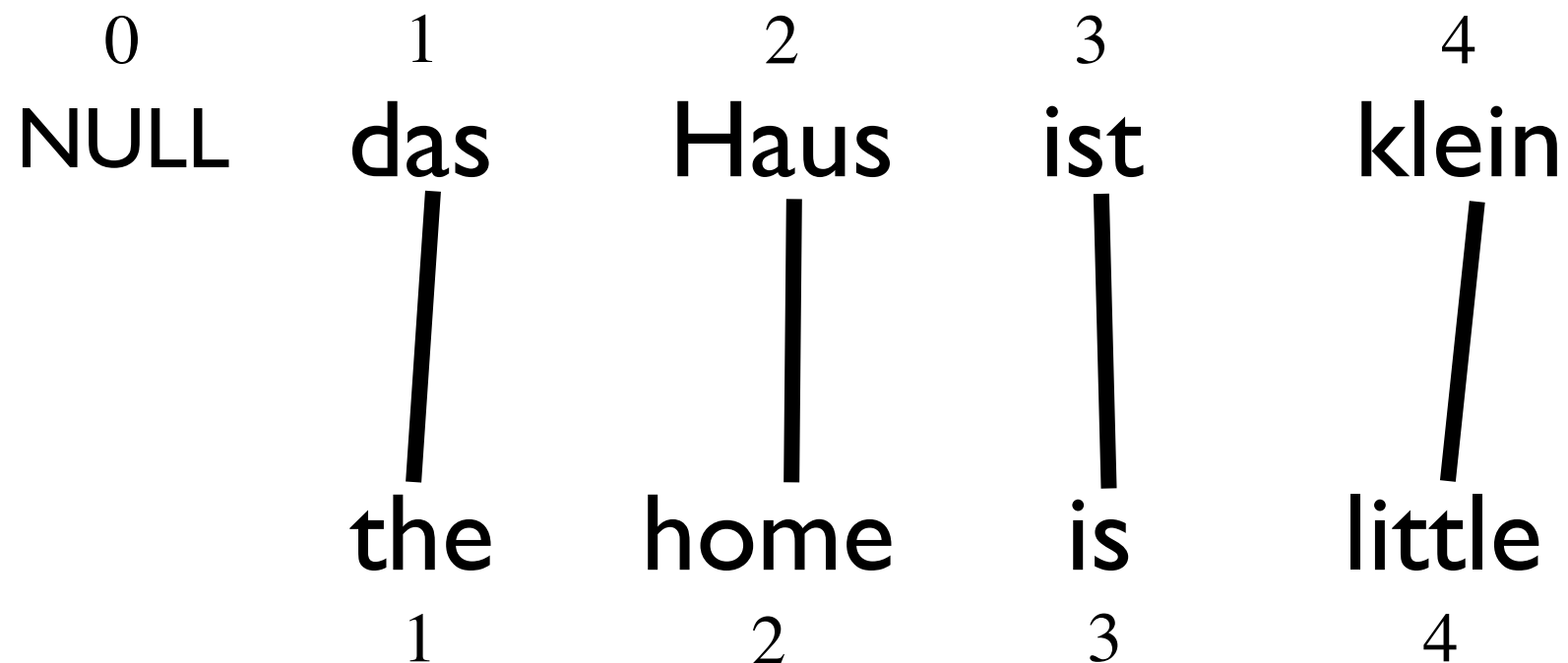
Finding the Viterbi Alignment



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

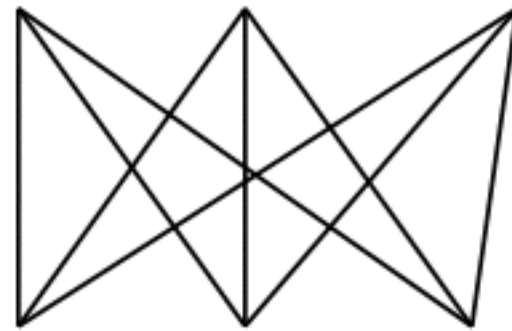
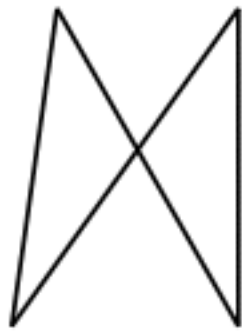
- pick some random (or uniform) parameters
- Repeat until you get bored (~ 5 iterations for lexical translation models)
 - using your current parameters, compute “expected” alignments for every target word token in the training data

$$p(a_i \mid \mathbf{e}, \mathbf{f}) \quad (\text{on board})$$

- 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 ...

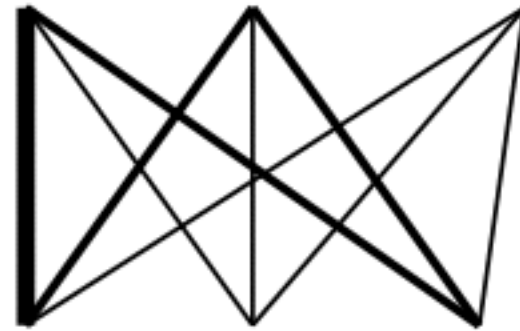


... the house ... the blue house ... the flower ...

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

EM for Model I

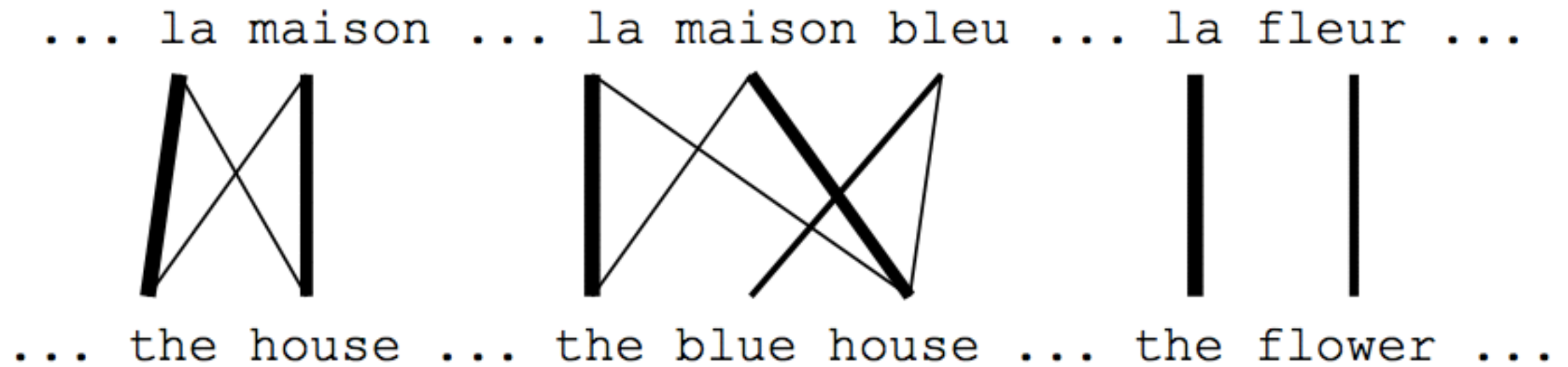
... la maison ... la maison blue ... la fleur ...



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

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

EM for Model 1



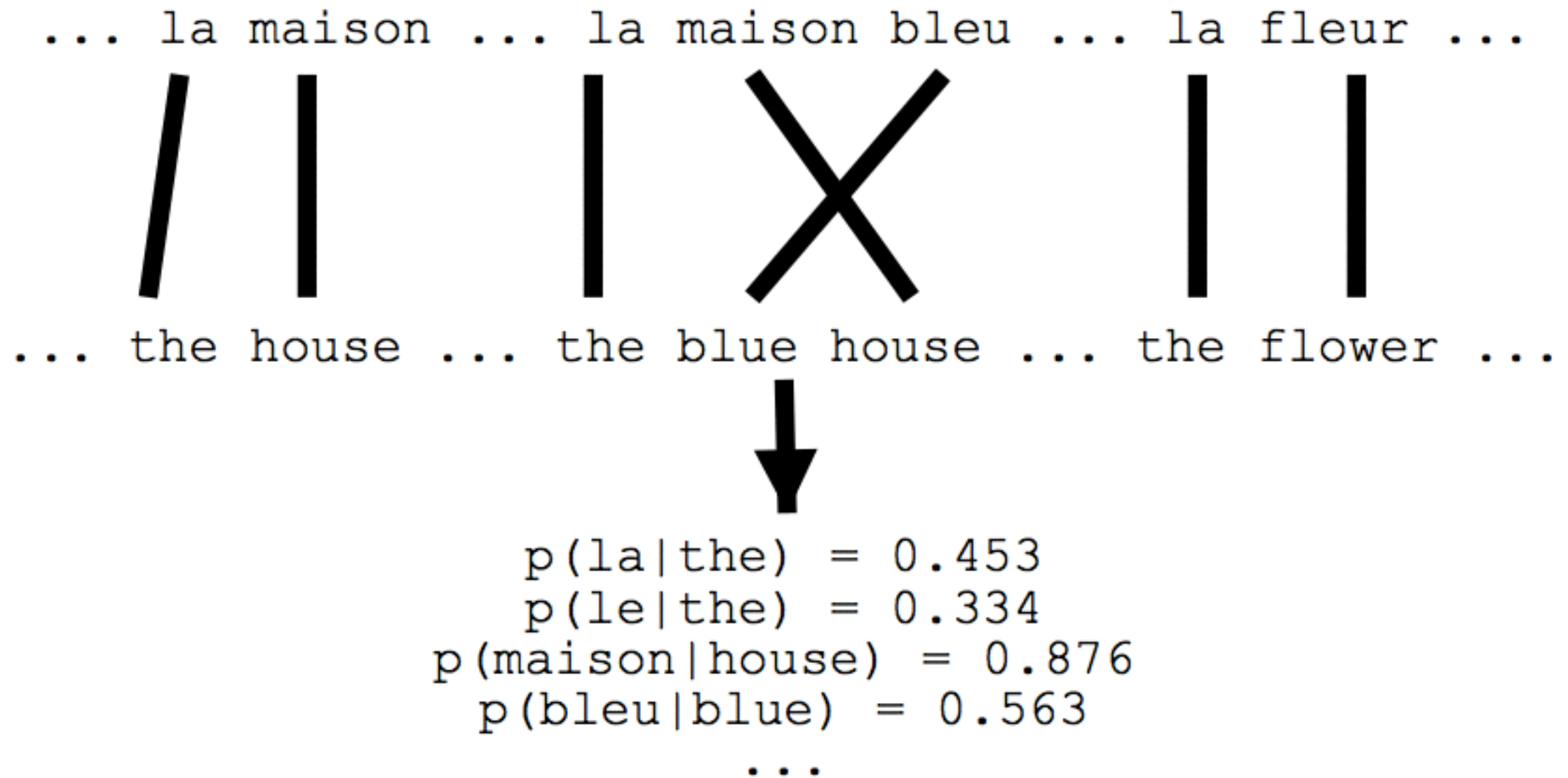
- After another iteration
- It becomes apparent that alignments, e.g., between **fleur** and **flower** are more likely (pigeon hole principle)

EM for Model 1

... la maison ... la maison bleu ... la fleur ...
/ | | X | |
... the house ... the blue house ... the flower ...

- Convergence
- Inherent hidden structure revealed by EM

EM for Model 1



- 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

Evaluation

- Since we have a probabilistic model, we can evaluate **perplexity**.

$$\text{PPL} = 2^{-\frac{1}{\sum_{(\mathbf{e}, \mathbf{f}) \in \mathcal{D}} |\mathbf{e}|} \log \prod_{(\mathbf{e}, \mathbf{f}) \in \mathcal{D}} p(\mathbf{e}|\mathbf{f})}$$

	Iter 1	Iter 2	Iter 3	Iter 4	...	Iter ∞
-log likelihood	-	7.66	7.21	6.84	...	-6
perplexity	-	2.42	2.3	2.21	...	2

Alignment Error Rate

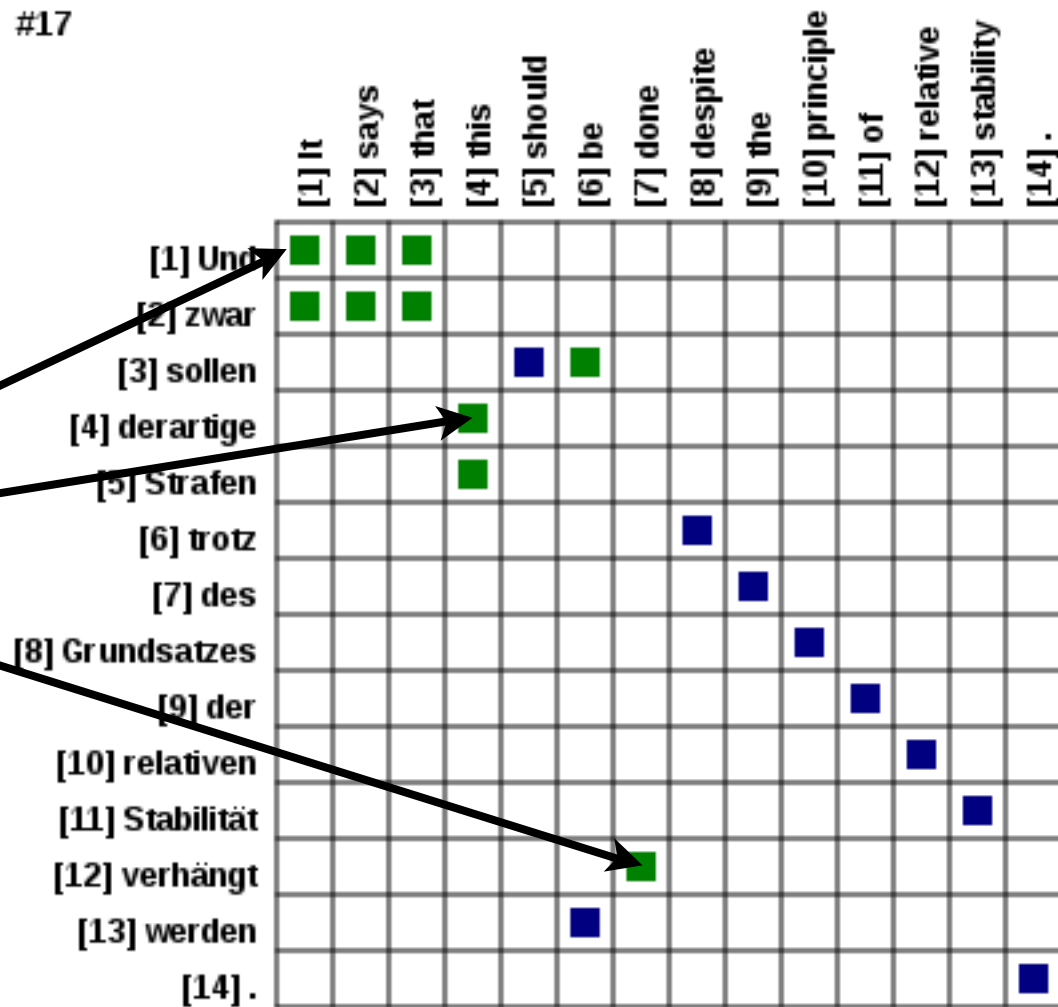
#17

[illegible]

Alignment Error Rate

#17

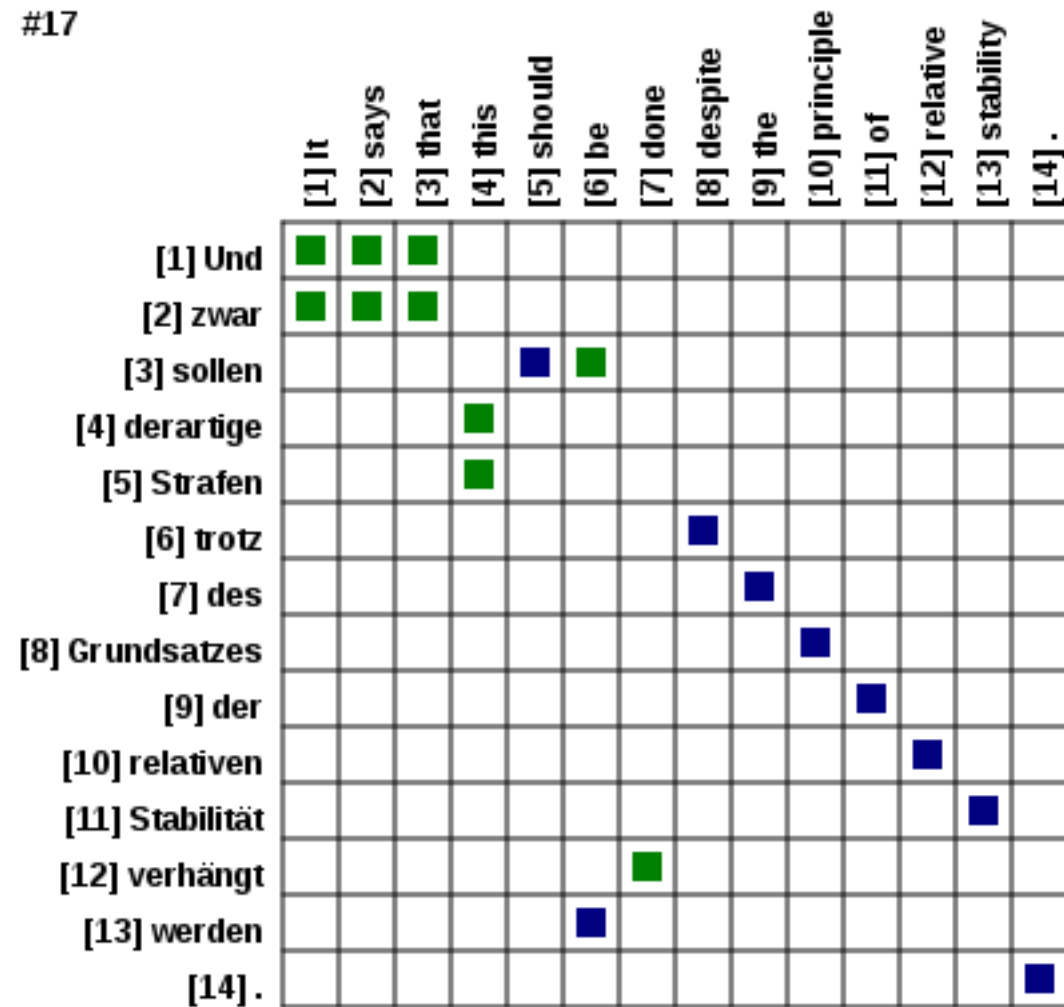
Possible links

$$P$$


Alignment Error Rate

#17

Possible links

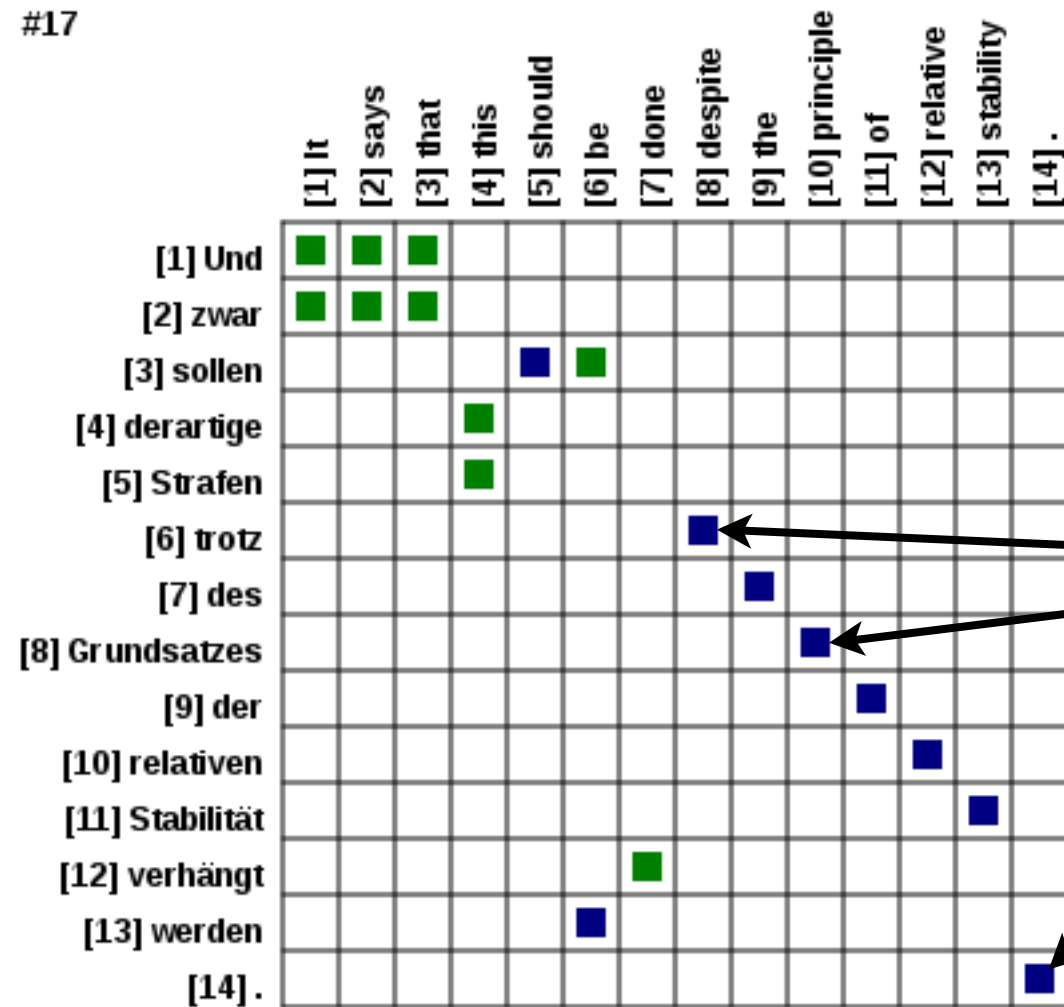
$$P$$


Alignment Error Rate

#17

Possible links

P



Sure links

S

Alignment Error Rate

#17

Possible links

$$P$$
[illegible]

Sure links

S

Alignment Error Rate

#17

Possible links

P

Sure links

S

	[1] It	[2] says	[3] that	[4] this	[5] should	[6] be	[7] done	[8] despite	[9] the	[10] principle	[11] of	[12] relative	[13] stability	[14].
[1] Und	■	■	■											
[2] zwar	■	■	■											
[3] sollen					■	■								
[4] derartige				■										
[5] Strafen				■										
[6] trotz								■						
[7] des									■					
[8] Grundsatzes										■				
[9] der											■			
[10] relativen												■		
[11] Stabilität													■	
[12] verhängt							■							
[13] werden						■								
[14].														■

$$\text{Precision}(A, P) = \frac{|P \cap A|}{|A|}$$

Alignment Error Rate

#17

Possible links

P

Sure links

S

	[1] It	[2] says	[3] that	[4] this	[5] should	[6] be	[7] done	[8] despite	[9] the	[10] principle	[11] of	[12] relative	[13] stability	[14].
[1] Und	■	■	■											
[2] zwar	■	■	■											
[3] sollen					■	■								
[4] derartige				■										
[5] Strafen				■										
[6] trotz								■						
[7] des									■					
[8] Grundsatzes										■				
[9] der											■			
[10] relativen												■		
[11] Stabilität													■	
[12] verhängt							■							
[13] werden						■								
[14].														■

$$\text{Precision}(A, P) = \frac{|P \cap A|}{|A|}$$

$$\text{Recall}(A, S) = \frac{|S \cap A|}{|S|}$$

Alignment Error Rate

#17

Possible links

P

Sure links

S

	[1] it	[2] says	[3] that	[4] this	[5] should	[6] be	[7] done	[8] despite	[9] the	[10] principle	[11] of	[12] relative	[13] stability	[14].
[1] Und	■	■	■											
[2] zwar	■	■	■											
[3] sollen					■	■								
[4] derartige				■										
[5] Strafen				■										
[6] trotz								■						
[7] des									■					
[8] Grundsatzes										■				
[9] der											■			
[10] relativen												■		
[11] Stabilität													■	
[12] verhängt							■							
[13] werden						■								
[14].														■

$$\text{Precision}(A, P) = \frac{|P \cap A|}{|A|}$$

$$\text{Recall}(A, S) = \frac{|S \cap A|}{|S|}$$

$$\text{AER}(A, P, S) = 1 - \frac{|S \cap A| + |P \cap A|}{|S| + |A|}$$

Alignment Error Rate

#17

Possible links

P

Sure links

S

	[1] it	[2] says	[3] that	[4] this	[5] should	[6] be	[7] done	[8] despite	[9] the	[10] principle	[11] of	[12] relative	[13] stability	[14].
[1] Und	■	■	■											
[2] zwar	■	■	■											
[3] sollen					■	■								
[4] derartige				■										
[5] Strafen				■										
[6] trotz								■						
[7] des									■					
[8] Grundsatzes										■				
[9] der											■			
[10] relativen												■		
[11] Stabilität													■	
[12] verhängt							■							
[13] werden						■								
[14].														■

$$\text{Precision}(A, P) = \frac{|P \cap A|}{|A|}$$

$$\text{Recall}(A, S) = \frac{|S \cap A|}{|S|}$$

$$\text{AER}(A, P, S) = 1 - \frac{|S \cap A| + |P \cap A|}{|S| + |A|}$$

Alignment Error Rate

#17

Possible links

P

Sure links

S

	[1] It	[2] says	[3] that	[4] this	[5] should	[6] be	[7] done	[8] despite	[9] the	[10] principle	[11] of	[12] relative	[13] stability	[14].
[1] Und	■	■	■											
[2] zwar	■	■	■											
[3] sollen					■	■								
[4] derartige				■										
[5] Strafen				■										
[6] trotz								■						
[7] des									■					
[8] Grundsatzes										■				
[9] der											■			
[10] relativen												■		
[11] Stabilität													■	
[12] verhängt							■							
[13] werden						■								
[14].														■

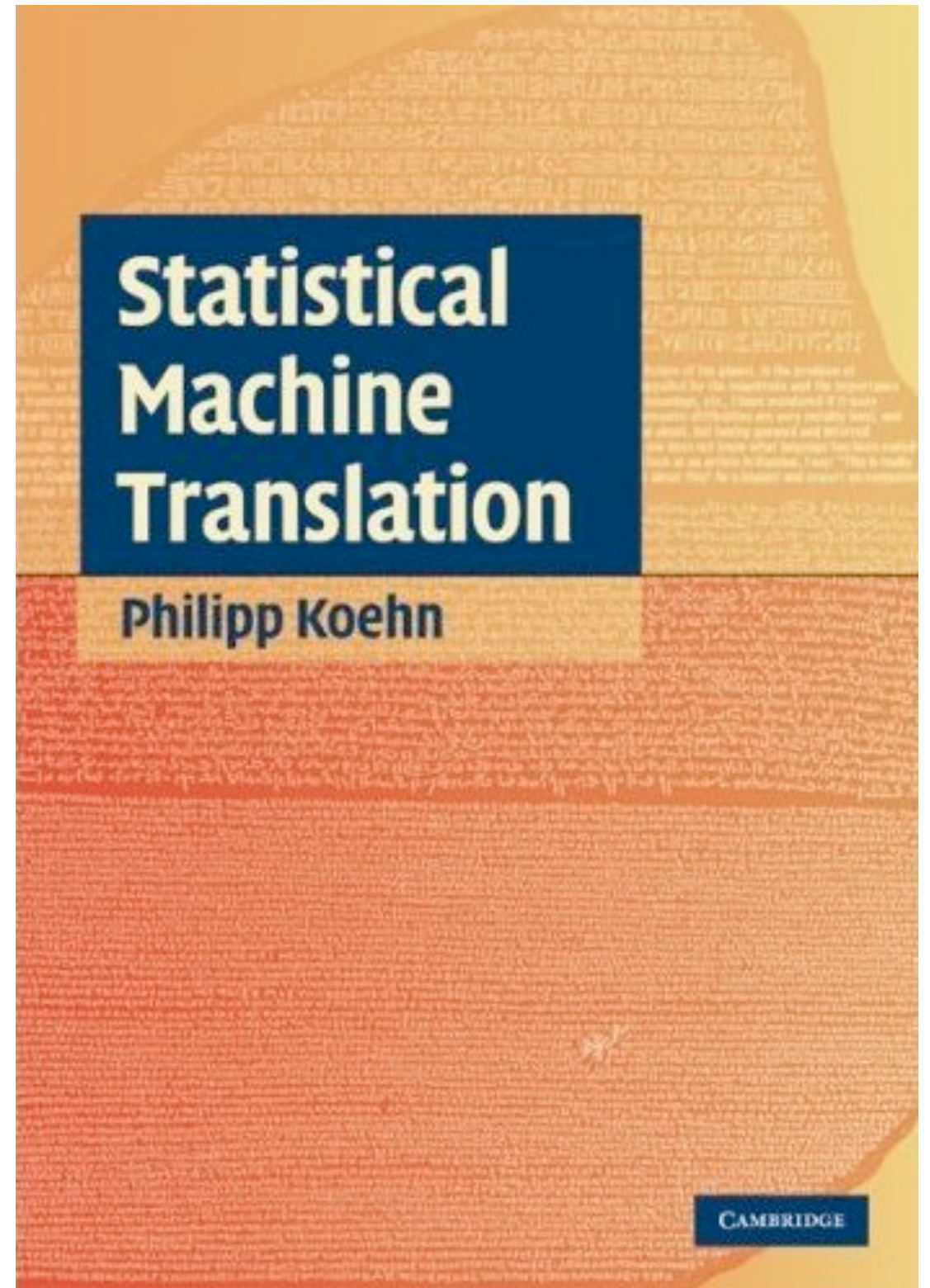
$$\text{Precision}(A, P) = \frac{|P \cap A|}{|A|}$$

$$\text{Recall}(A, S) = \frac{|S \cap A|}{|S|}$$

$$\text{AER}(A, P, S) = 1 - \frac{|S \cap A| + |P \cap A|}{|S| + |A|}$$

Reading

- Read Chapter 4 from the textbook (today we covered 4.1 and 4.2)



Announcements

- HW 1 is due in 1 week