

EM with Features

Nov. 19, 2013

Word Alignment

das Haus

ein Buch

das Buch

the house

a book

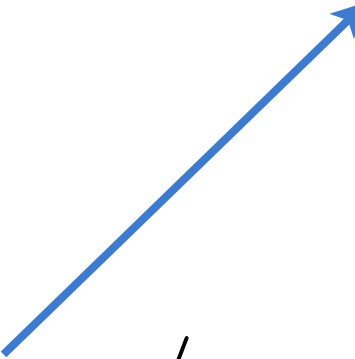
the book

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

- 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

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

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} .

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

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

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

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

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 alignments 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 alignments 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 alignments 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 alignments 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

<u>1</u>	<u>2</u>	<u>3</u>	<u>4</u>
----------	----------	----------	----------

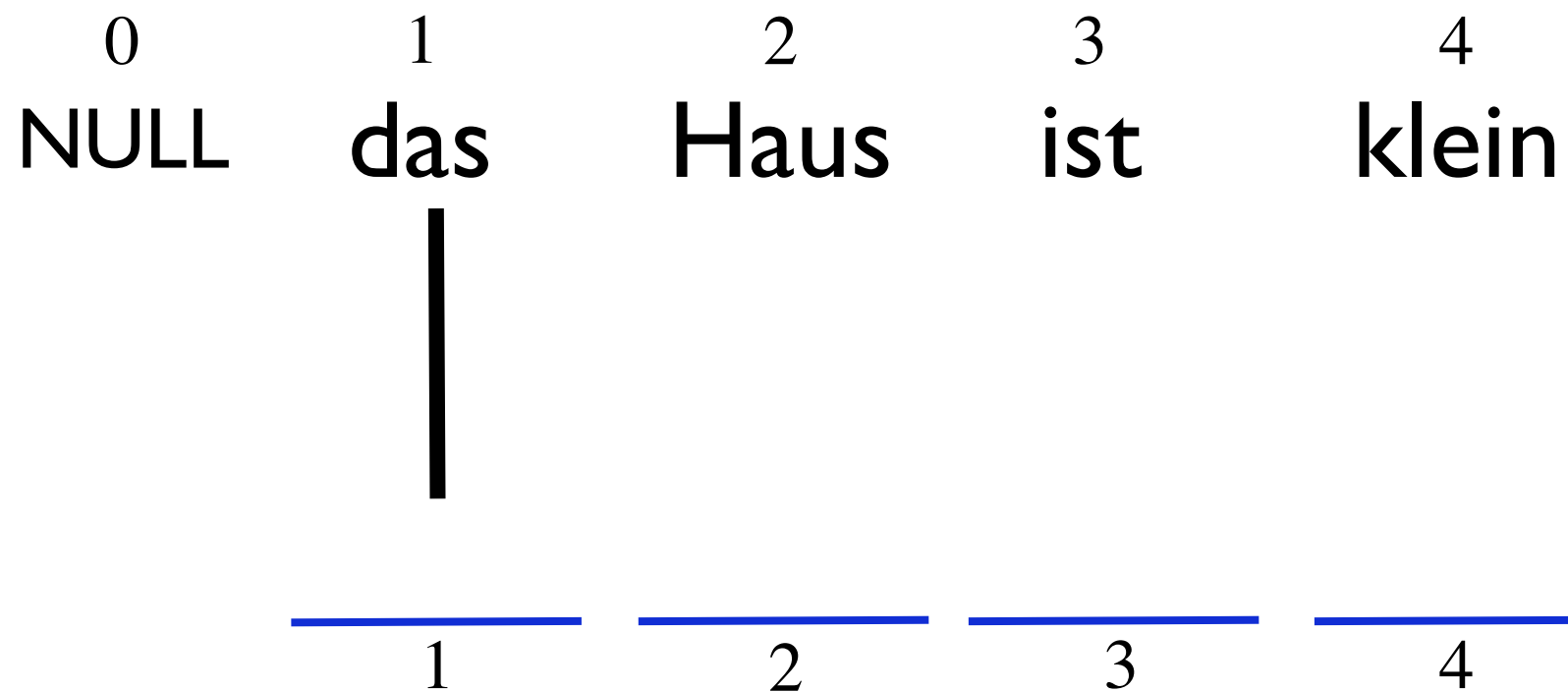
Start with a foreign sentence and a target length.

Example

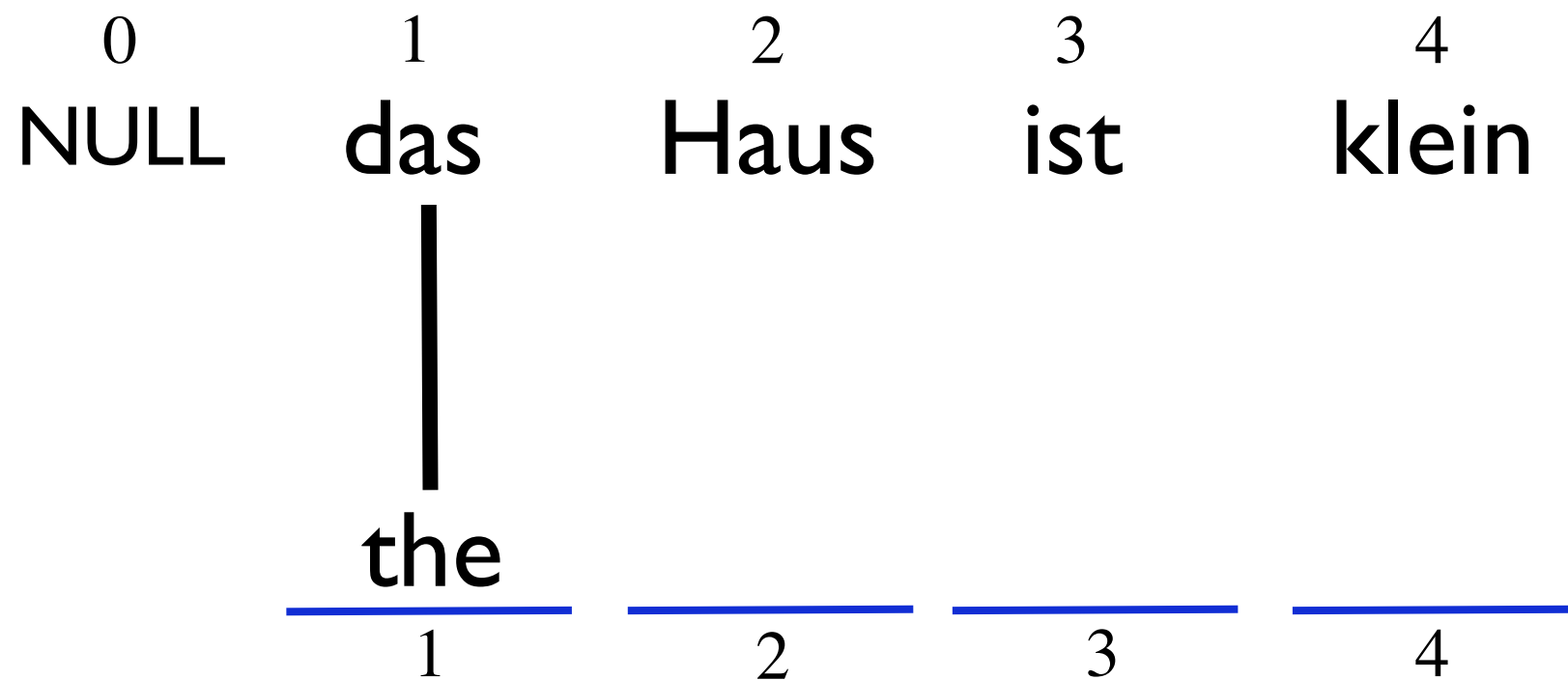
0	1	2	3	4
NULL	das	Haus	ist	klein

<u>1</u>	<u>2</u>	<u>3</u>	<u>4</u>
----------	----------	----------	----------

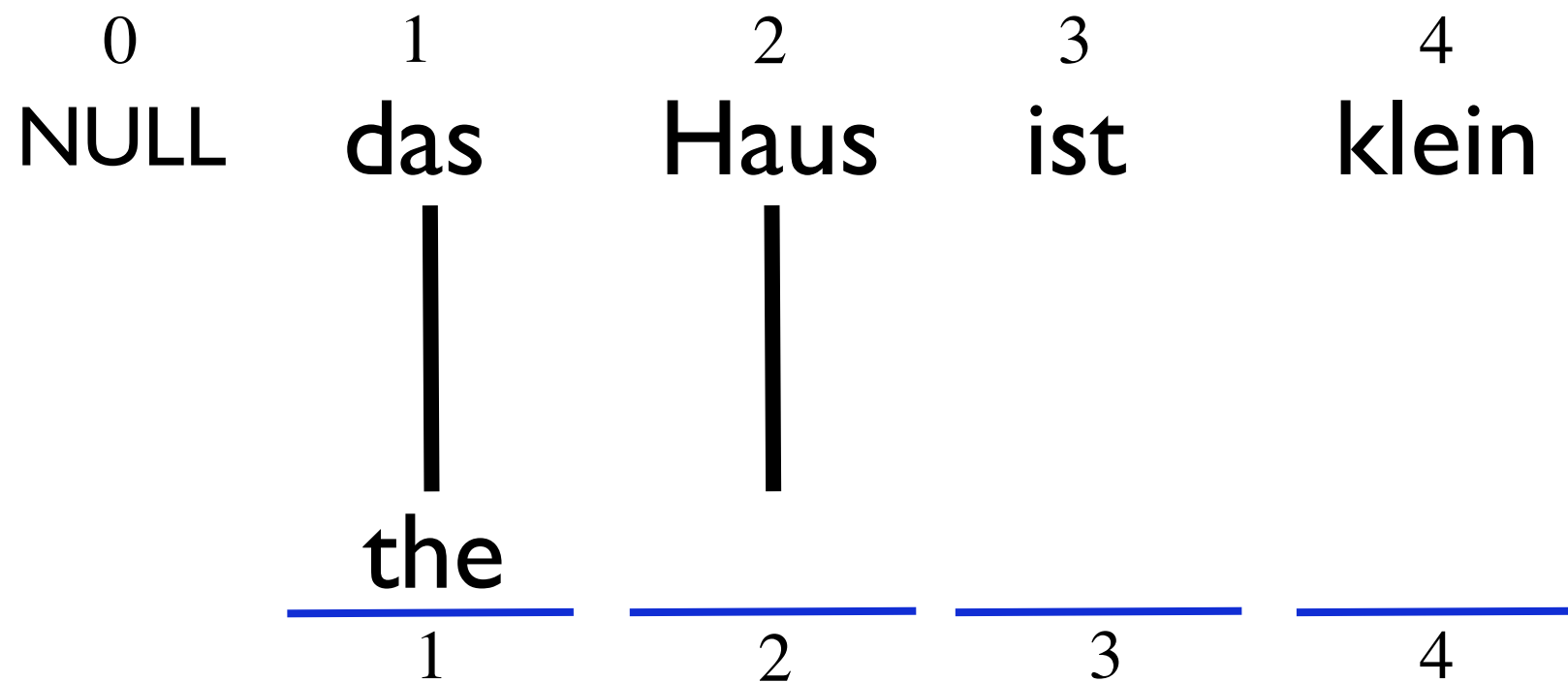
Example



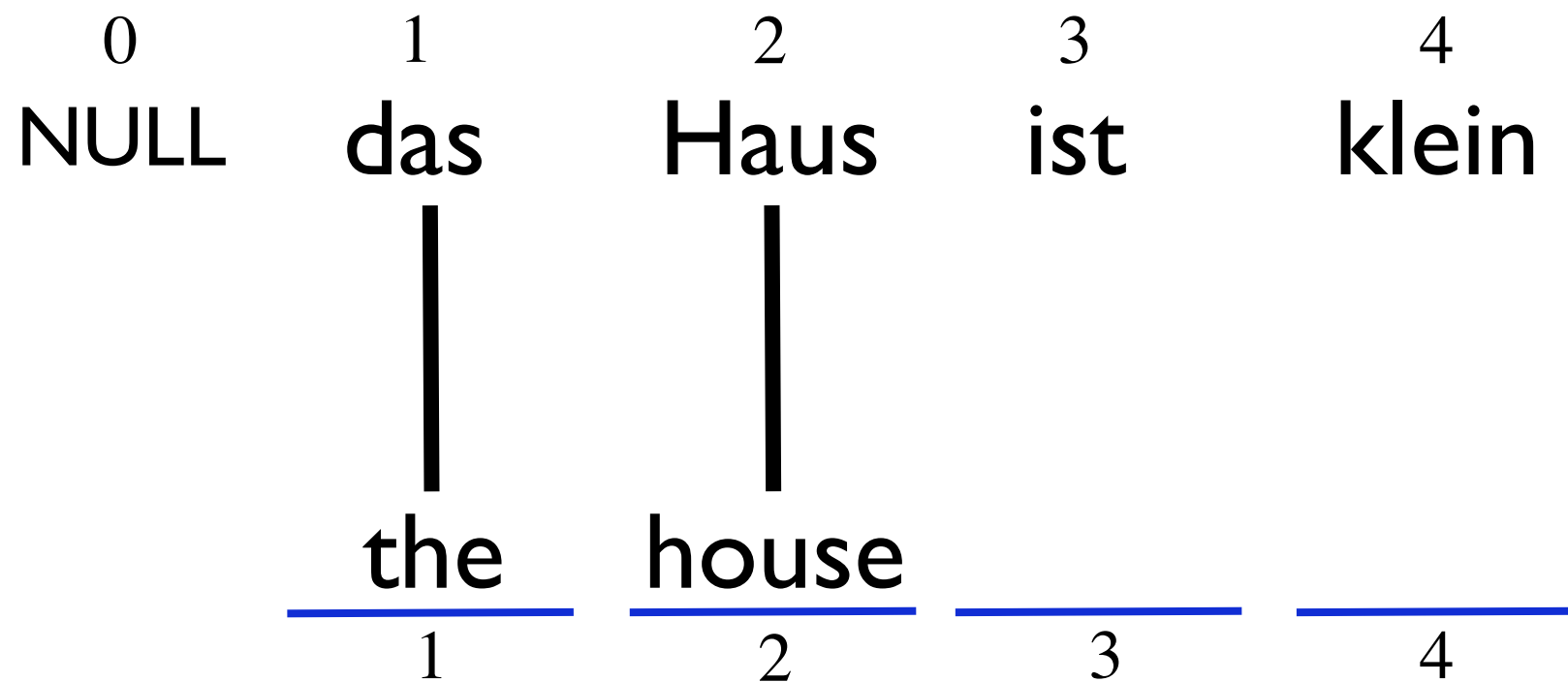
Example



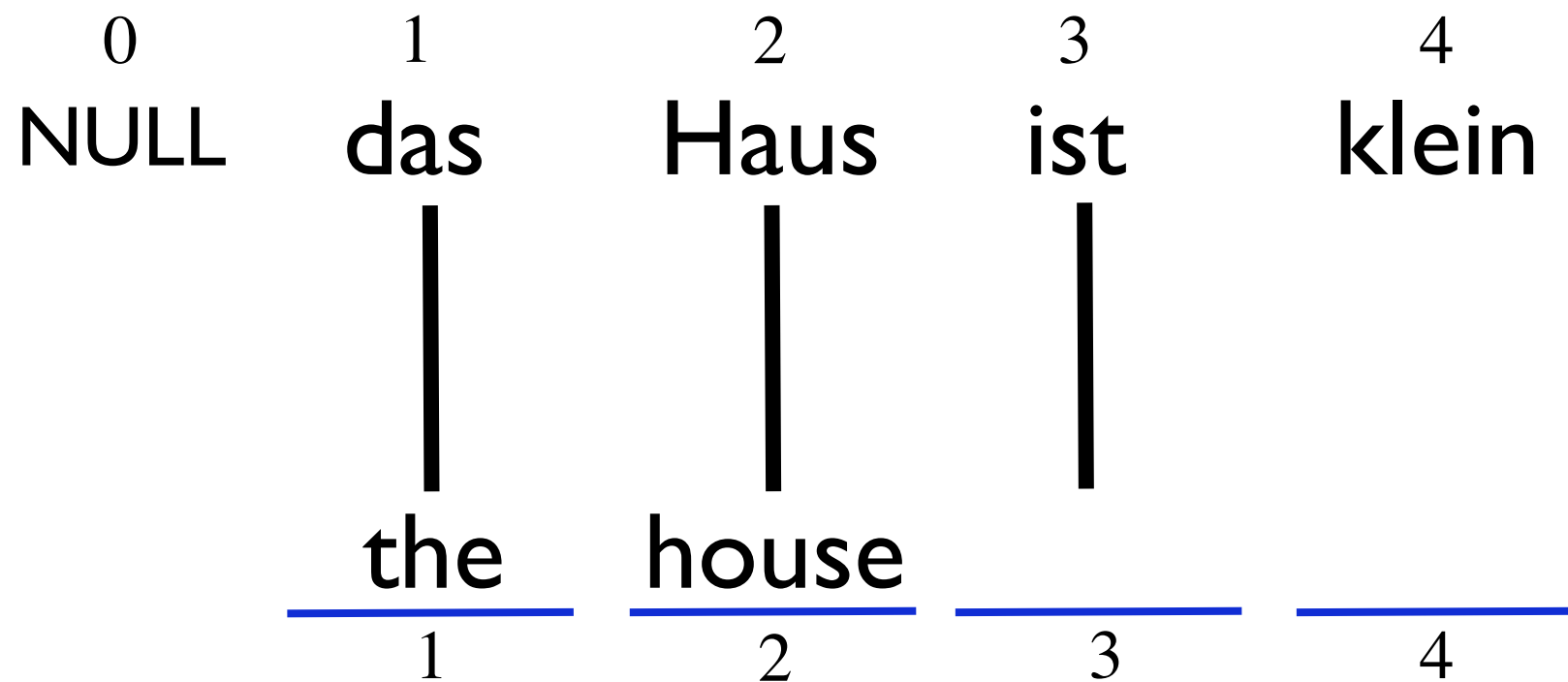
Example



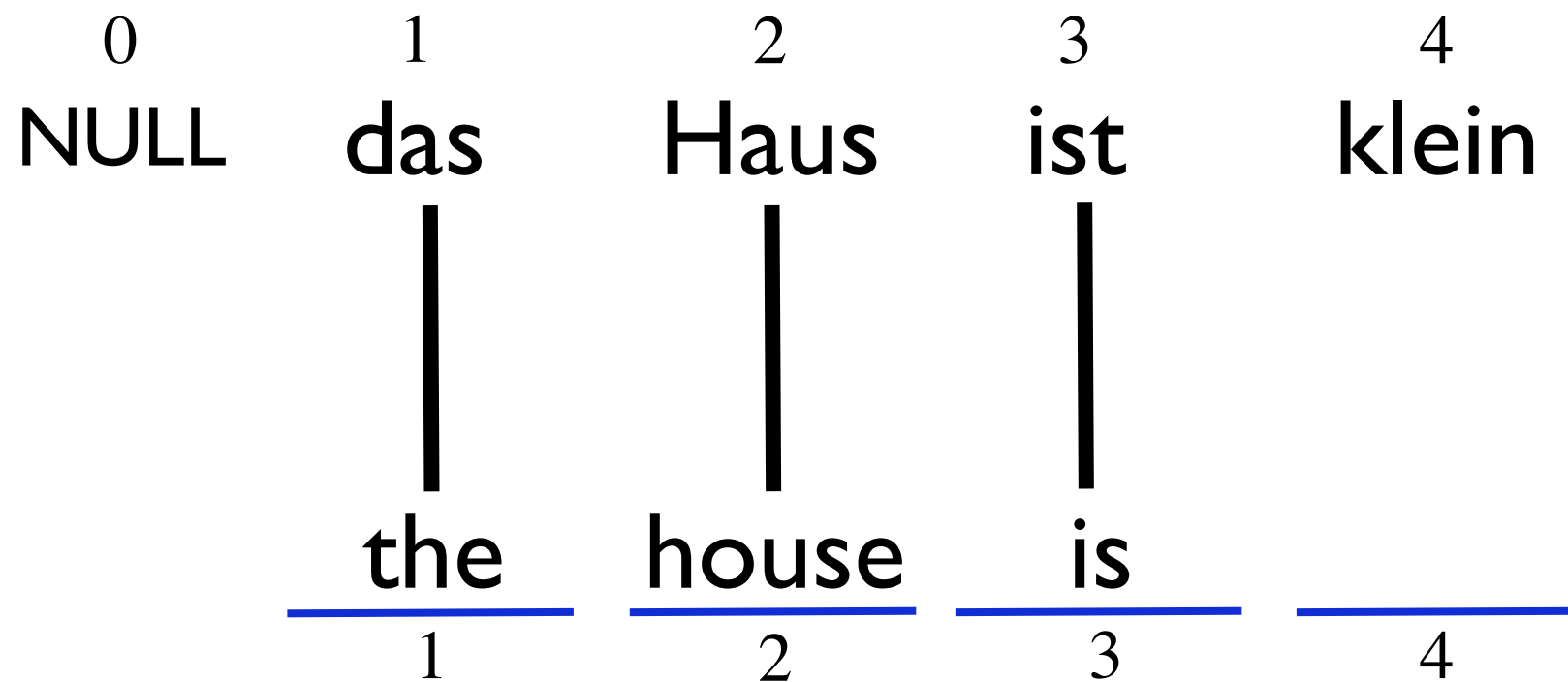
Example



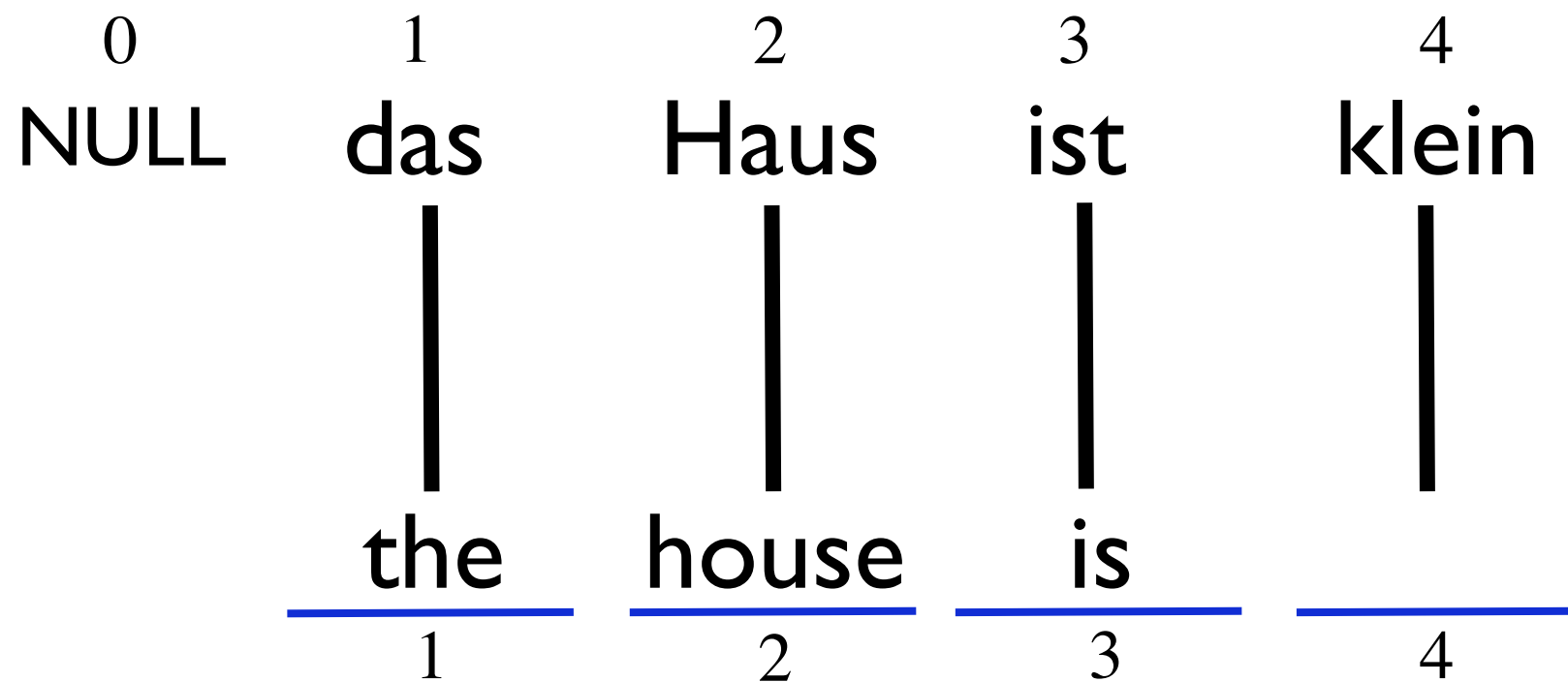
Example



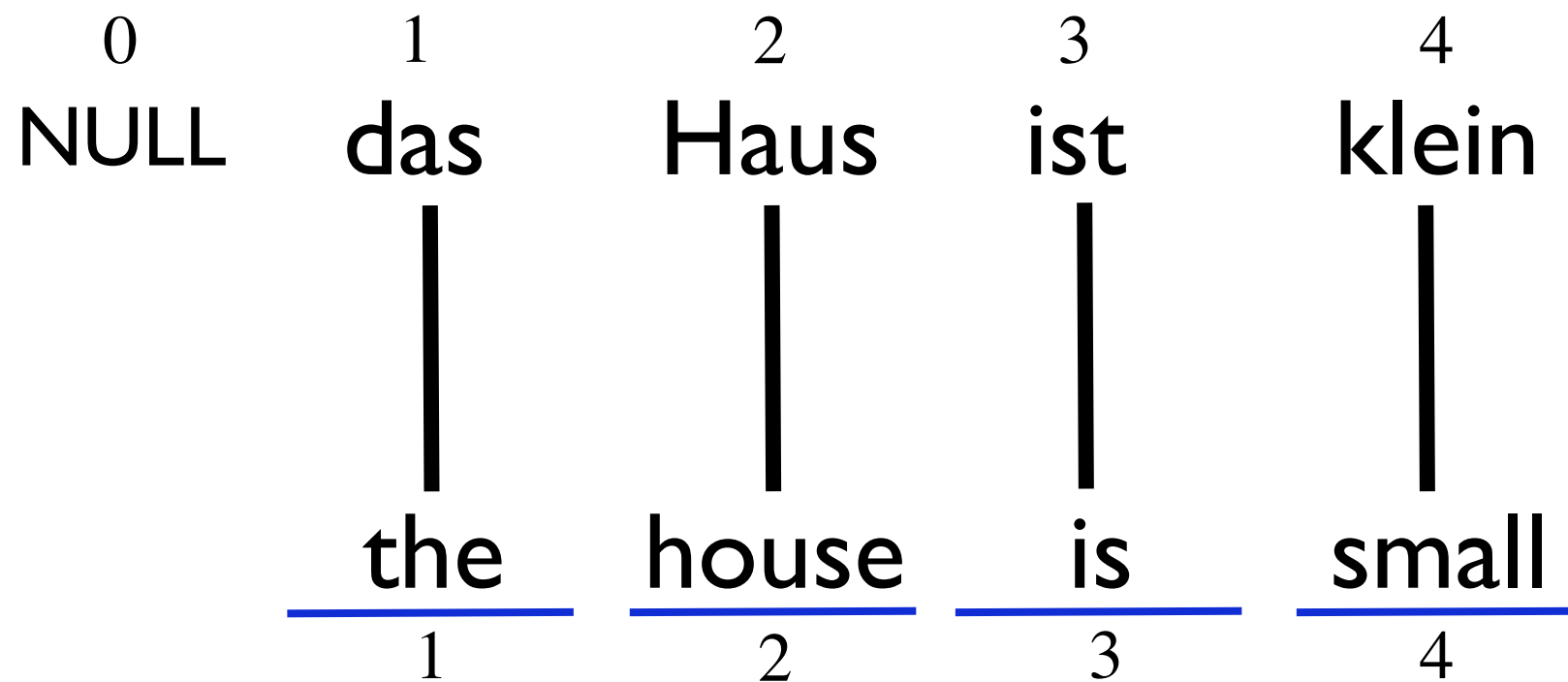
Example



Example



Example

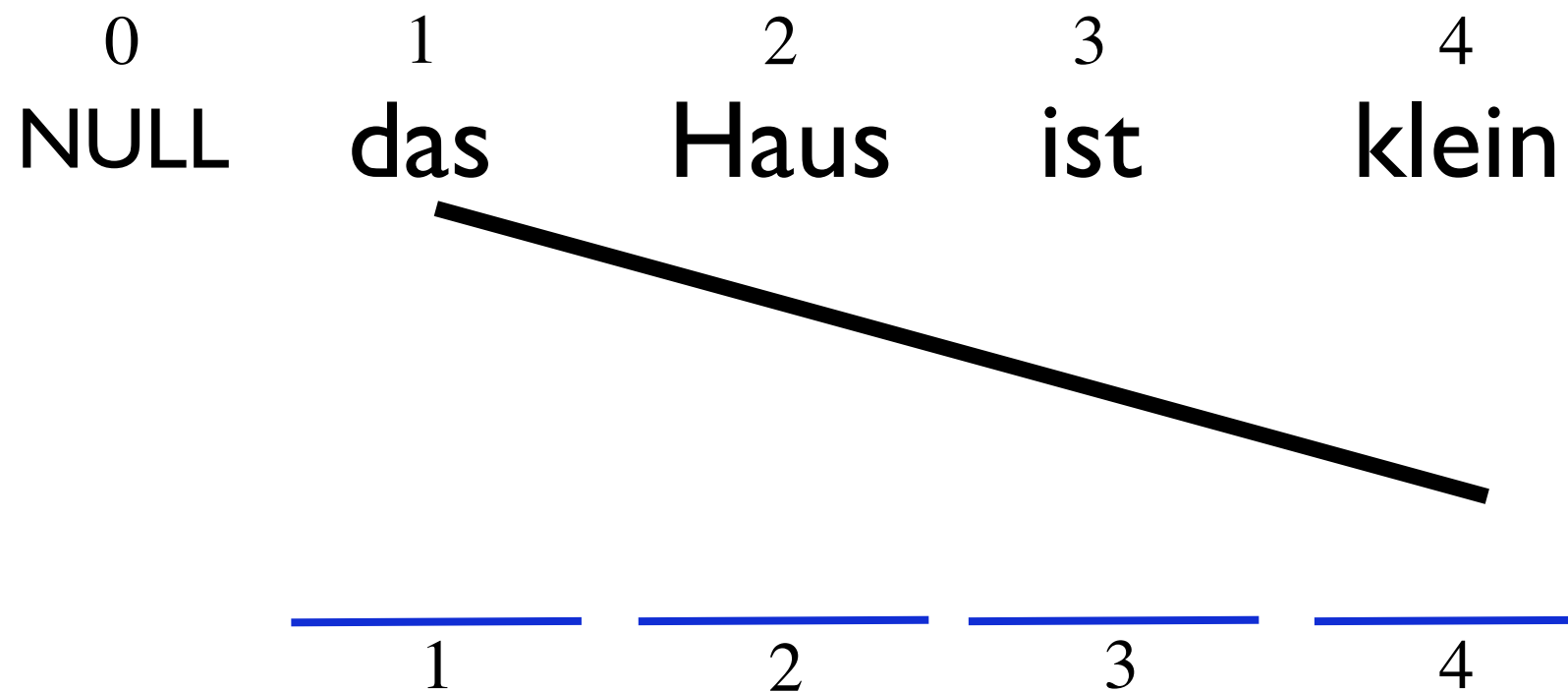


Example

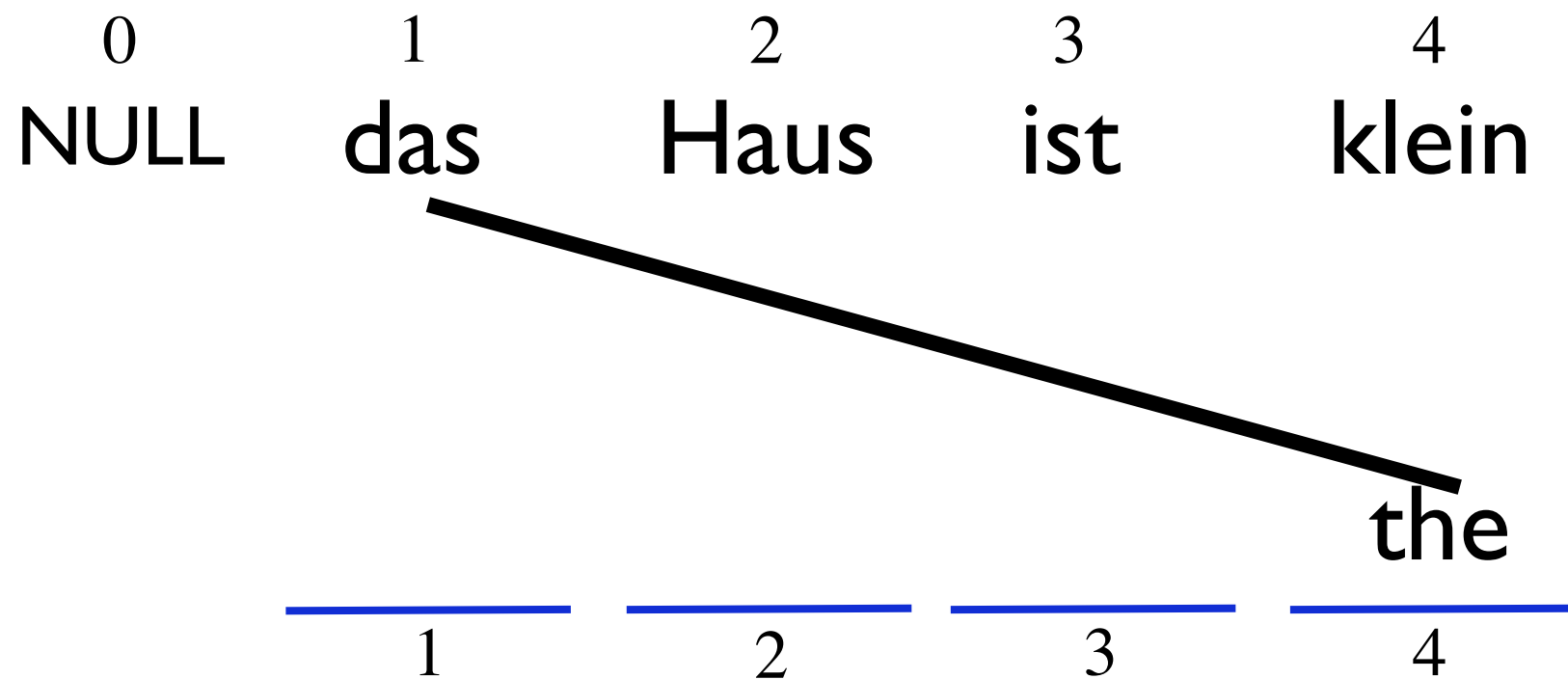
0	1	2	3	4
NULL	das	Haus	ist	klein

<u>1</u>	<u>2</u>	<u>3</u>	<u>4</u>
----------	----------	----------	----------

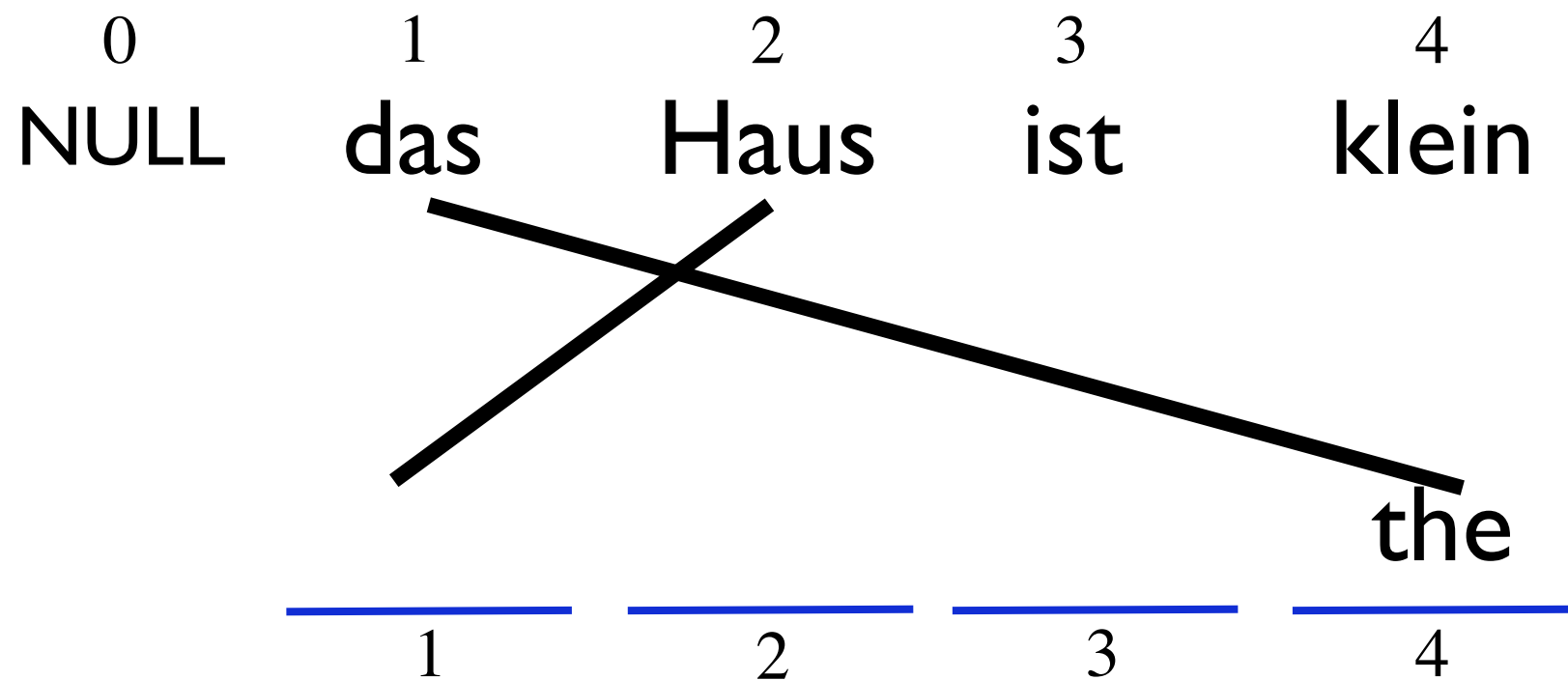
Example



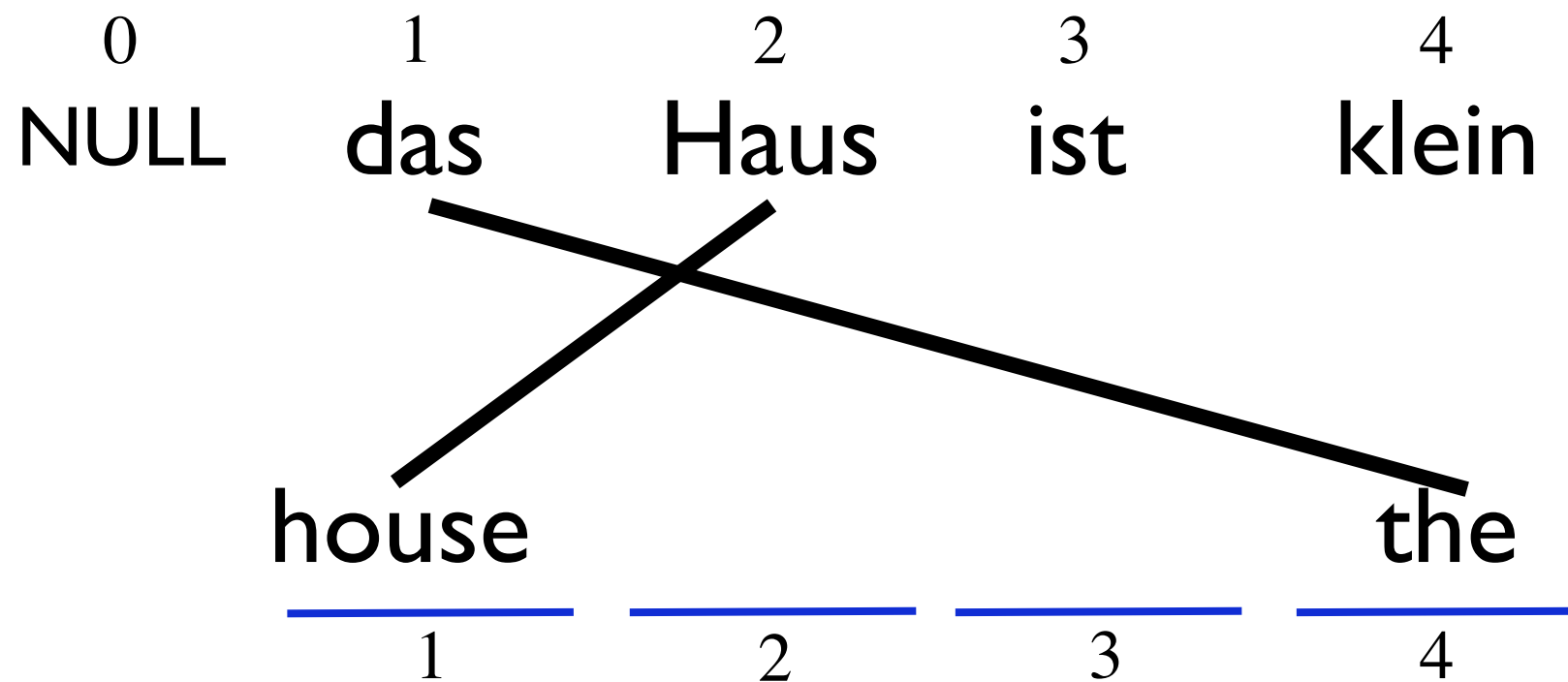
Example



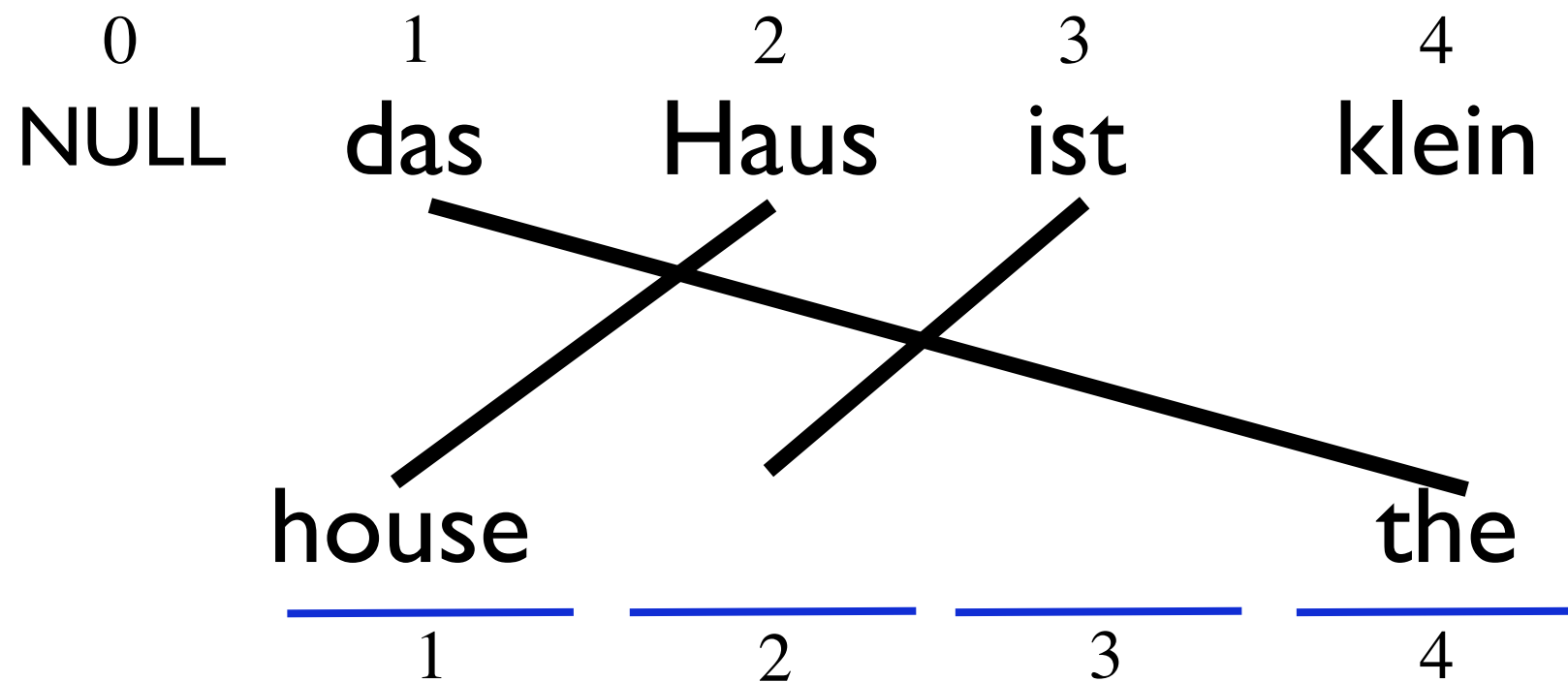
Example



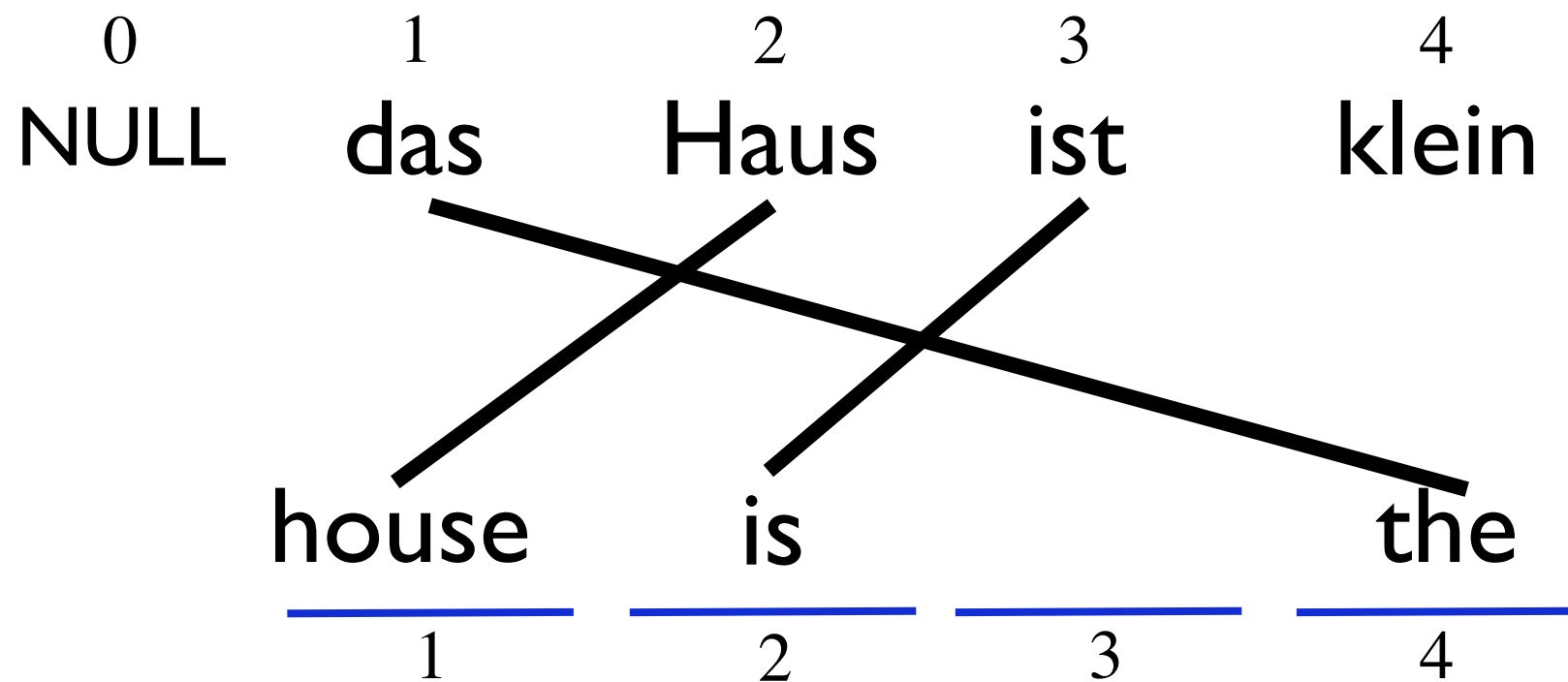
Example



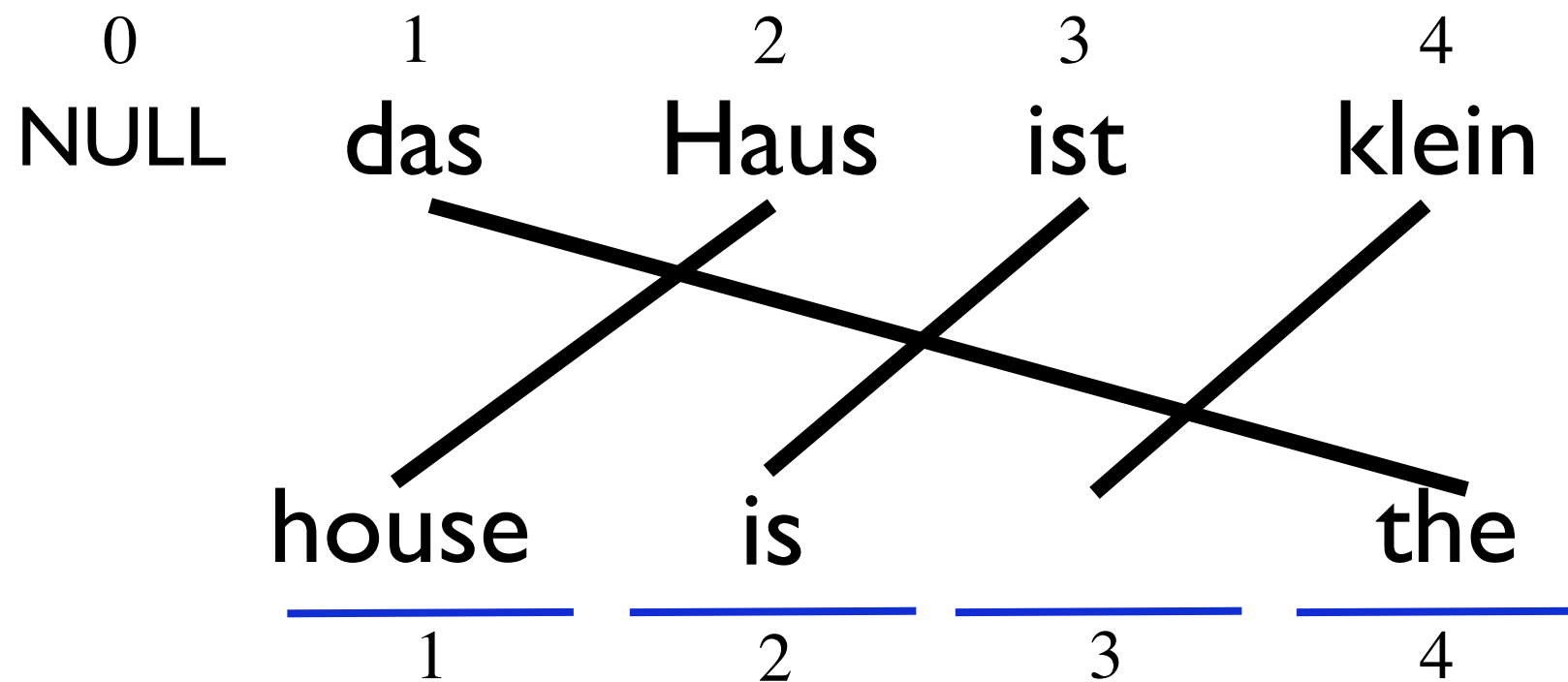
Example



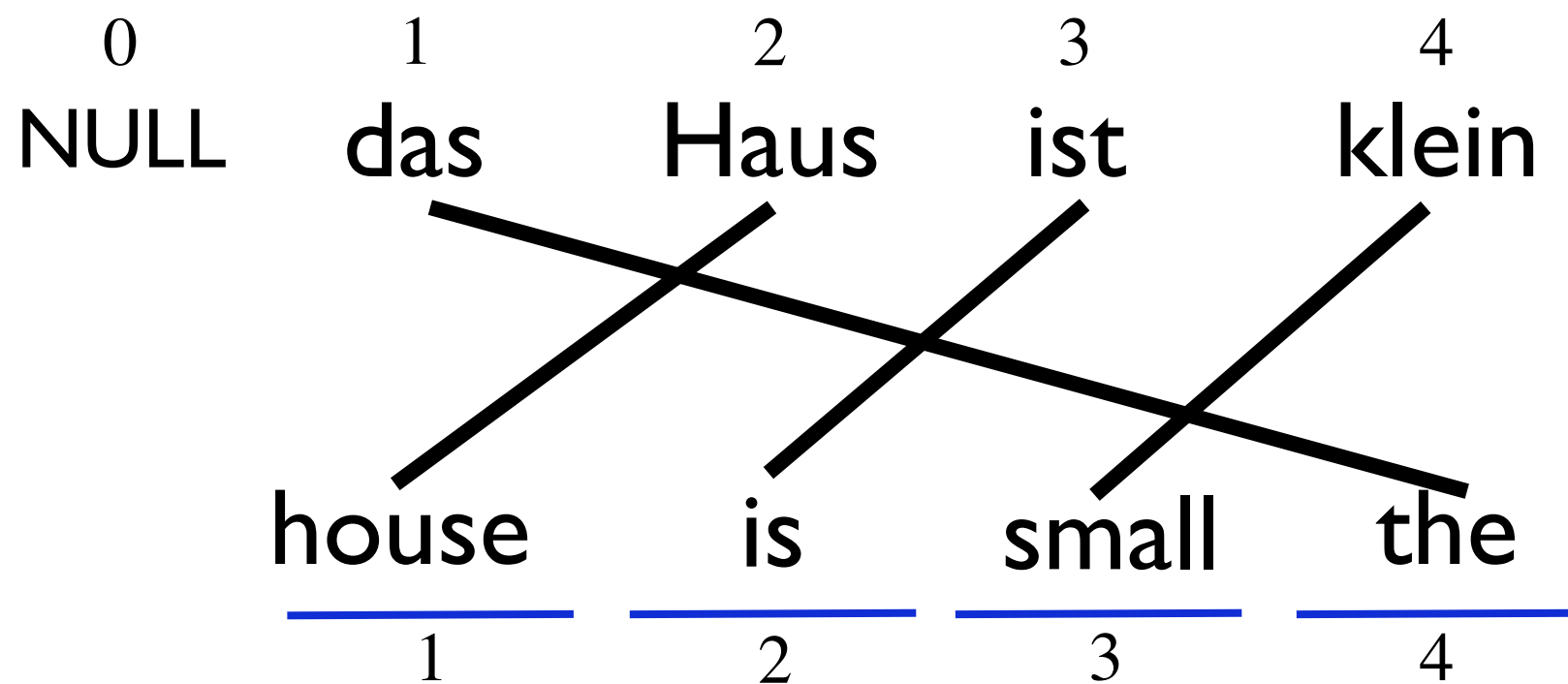
Example



Example



Example



das Haus
the house

das Buch
the book

ein Buch
a book

e	f	initial
the	das	0.25
book	das	0.25
house	das	0.25
the	buch	0.25
book	buch	0.25
a	buch	0.25
book	ein	0.25
a	ein	0.25
the	haus	0.25
house	haus	0.25

$freq(\text{Buch}, \text{book}) = ?$

$freq(\text{das}, \text{book}) = ?$

$freq(\text{ein}, \text{book}) = ?$

$freq(\text{Buch}, \text{book}) =$

$$\sum_i \mathbb{I}(\tilde{e}_i = \text{book}, \tilde{f}_{a_i} = \text{Buch})$$

$$\mathbb{E}_{p_{\mathbf{w}}(1)}(\mathbf{a} | \mathbf{f} = \text{das Buch}, \mathbf{e} = \text{the book}) \sum_i \mathbb{I}[e_i = \text{book}, f_{a_i} = \text{Buch}]$$

Convergence

das Haus
the house

das Buch
the book

ein Buch
a book

<i>e</i>	<i>f</i>	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	2.30	2.21	...	2

Hidden Markov Models

- GMMs, the aggregate bigram model, and Model I don't have conditional dependencies between random variables
- Let's consider an example of a model where this is not the case

$$p(\mathbf{x}) = \sum_{\mathbf{y}' \in \mathcal{Y}_{\mathbf{x}}} \eta(y_{|\mathbf{x}|} \rightarrow \text{STOP}) \prod_{i=1}^{|\mathbf{x}|} \eta(y_{i-1} \rightarrow y_i) \times \gamma(y_i \downarrow x_i)$$

EM for HMMs

- What statistics are sufficient to determine the parameter values?

$\text{freq}(q \downarrow x)$ How often does q emit x ?

$\text{freq}(q \rightarrow r)$ How often does q transition to r ?

$\text{freq}(q)$ How often do we visit q ?

EM for HMMs

- What statistics are sufficient to determine the parameter values?

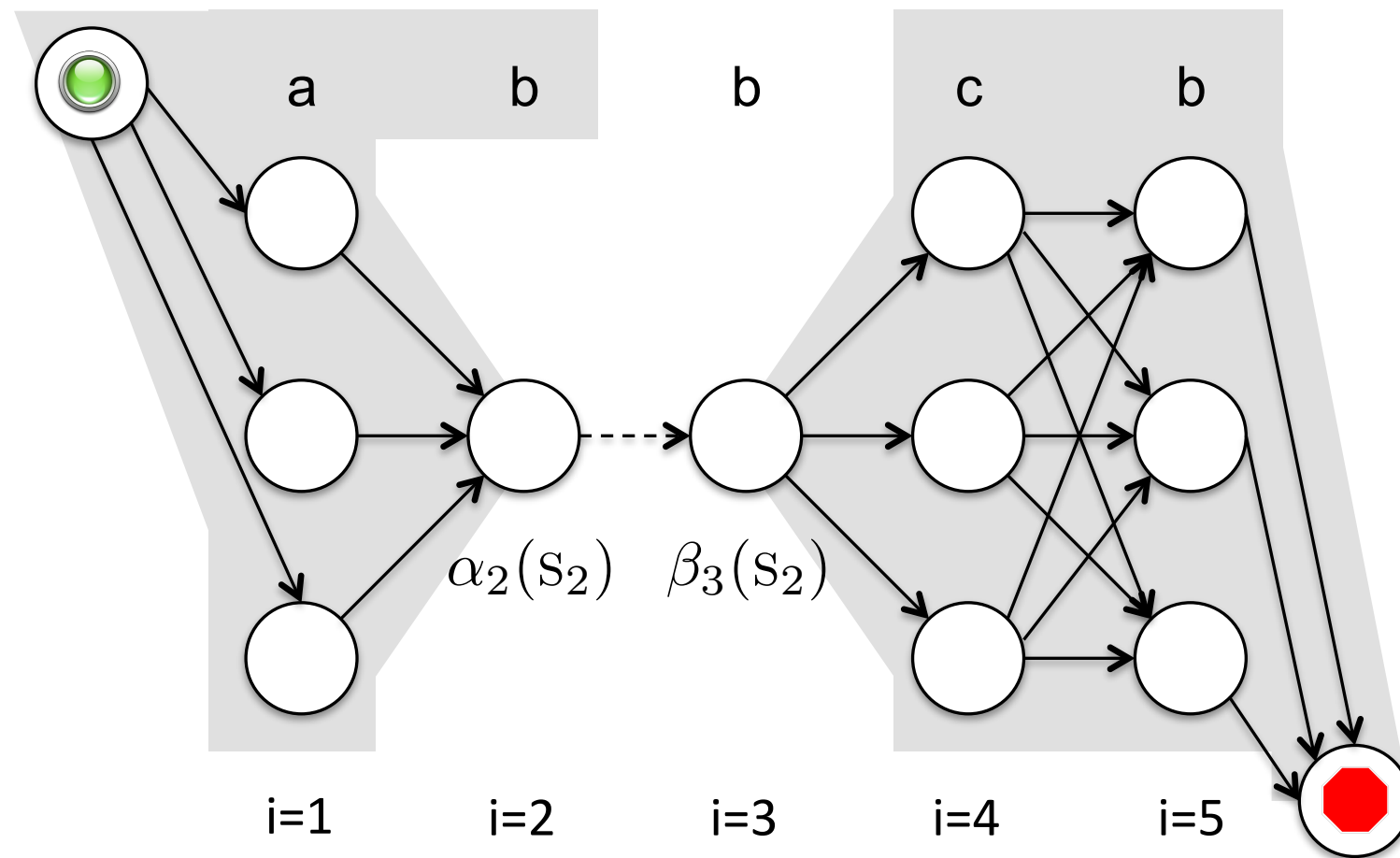
$\text{freq}(q \downarrow x)$ How often does q emit x ?

$\text{freq}(q \rightarrow r)$ How often does q transition to r ?

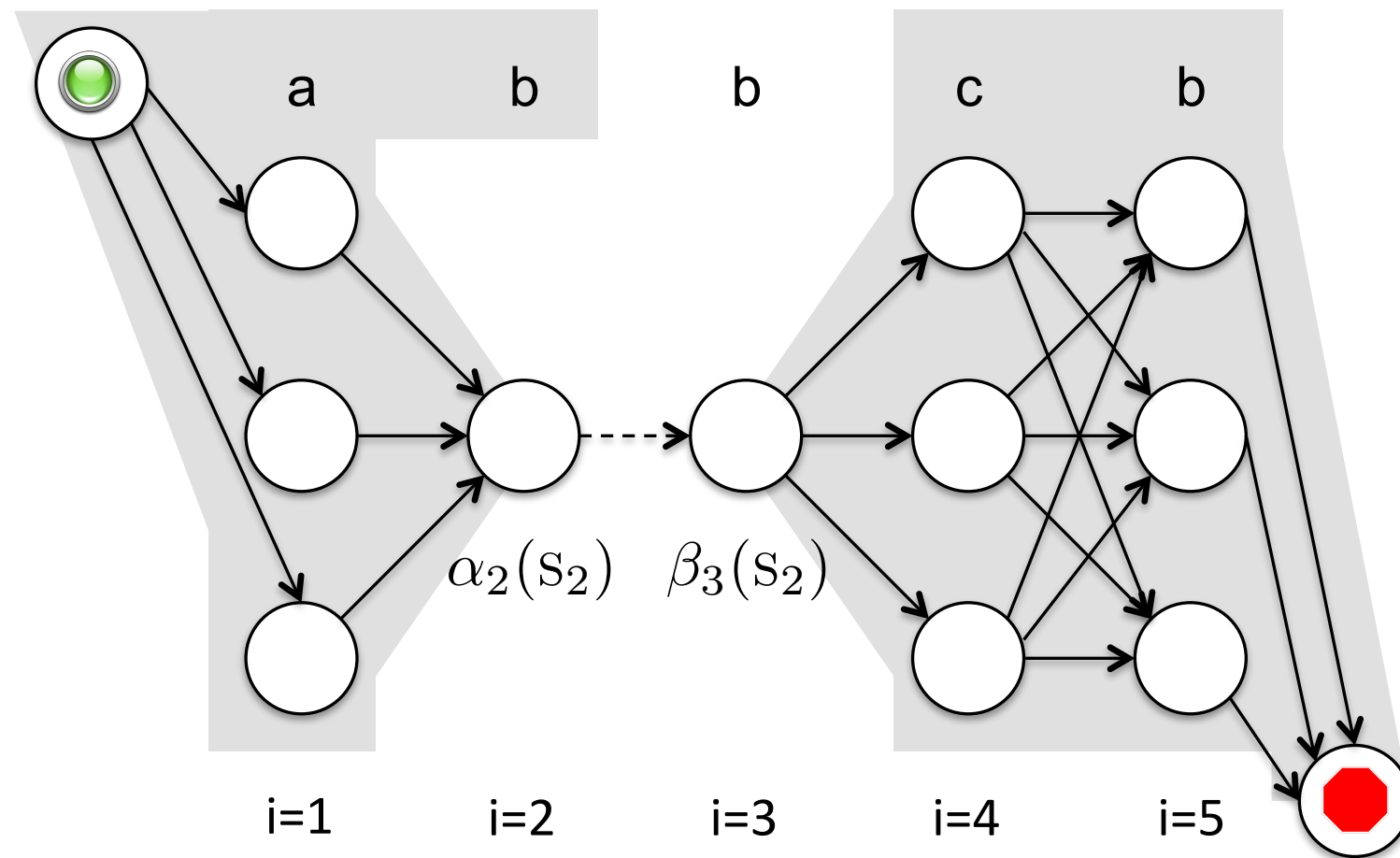
$\text{freq}(q)$ How often do we visit q ?

And of course...

$$\text{freq}(q) = \sum_{r \in Q} \text{freq}(q \rightarrow r)$$



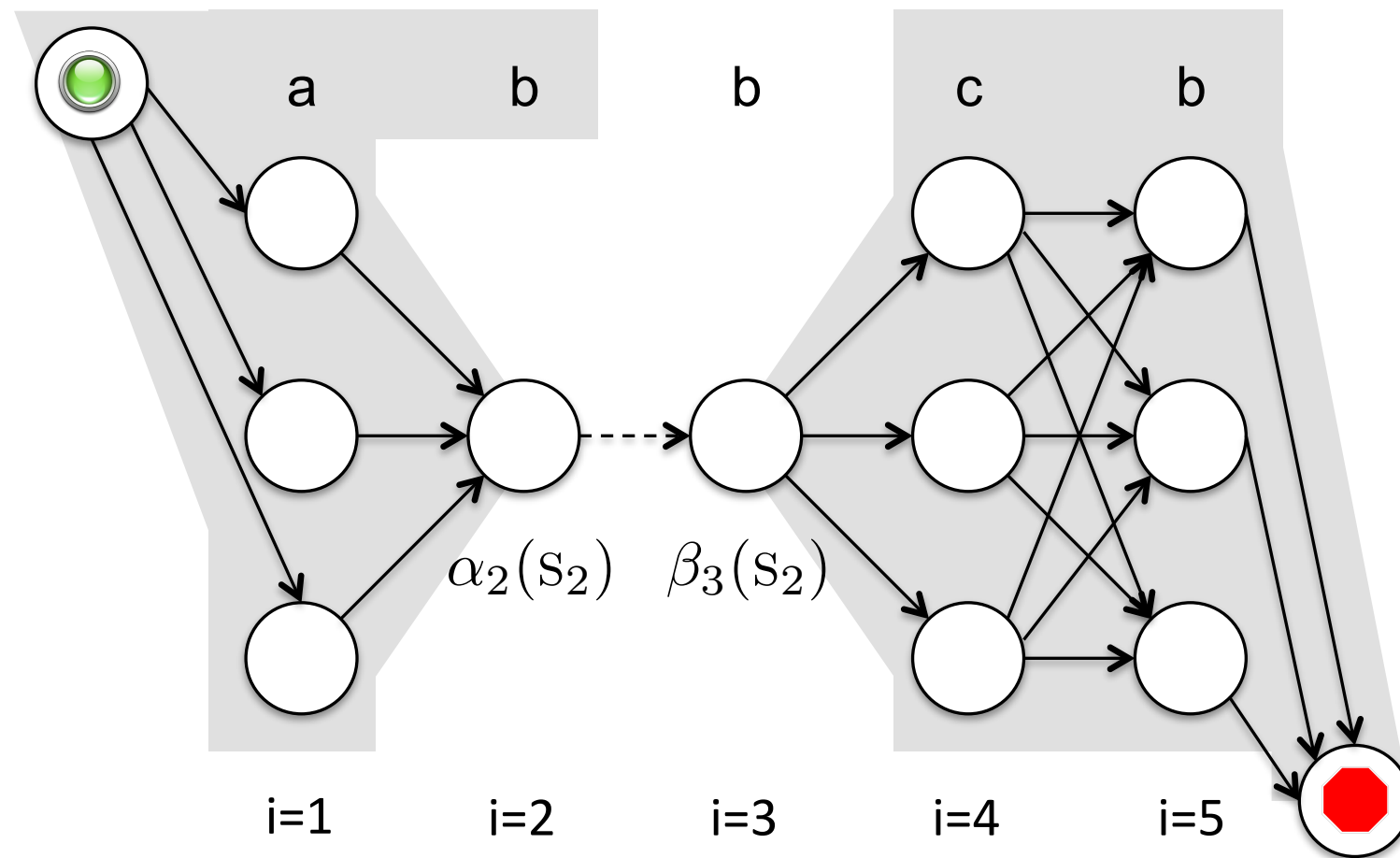
$$\begin{aligned}
 p(y_2 = q, y_3 = r \mid \mathbf{x}) &\propto p(y_2 = q, y_3 = r, \mathbf{x}) \\
 &= \frac{\alpha_2(q) \times \beta_3(r) \times \eta(q \rightarrow r) \times \eta(r \downarrow x_3)}{\sum_{q', r' \in \mathcal{Q}} \alpha_2(q') \times \beta_3(r') \times \eta(q' \rightarrow r') \times \eta(r' \downarrow x_3)} \\
 &= \frac{\alpha_2(q) \times \beta_3(r) \times \eta(q \rightarrow r) \times \eta(r \downarrow x_3)}{p(\mathbf{x}) = \alpha_{|\mathbf{x}|}(\text{STOP})}
 \end{aligned}$$



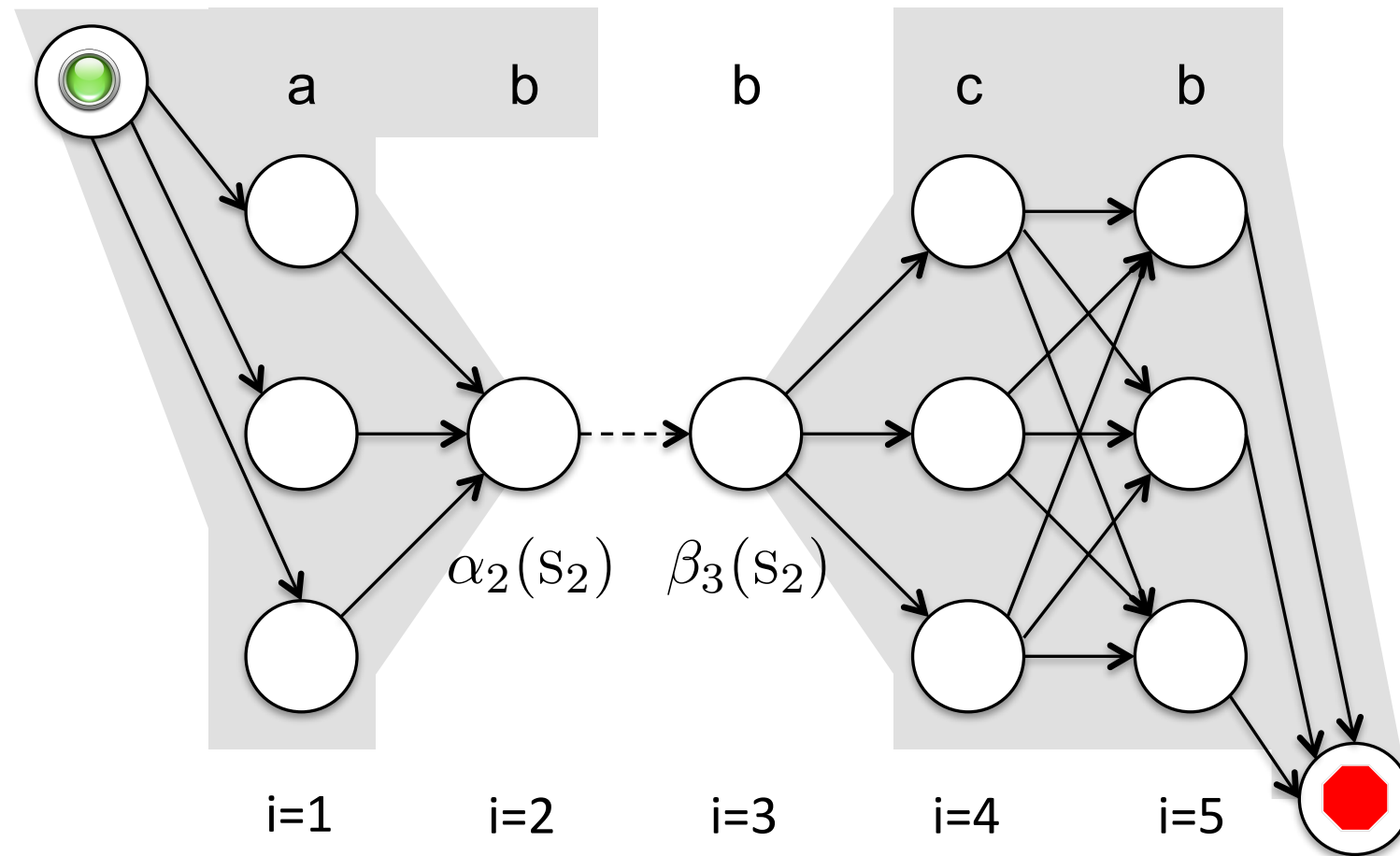
$$p(y_2 = q, y_3 = r \mid \mathbf{x}) \propto p(y_2 = q, y_3 = r, \mathbf{x})$$

$$= \frac{\alpha_2(q) \times \beta_3(r) \times \eta(q \rightarrow r) \times \eta(r \downarrow x_3)}{\sum_{q', r' \in \mathcal{Q}} \alpha_2(q') \times \beta_3(r') \times \eta(q' \rightarrow r') \times \eta(r' \downarrow x_3)}$$

$$= \frac{\alpha_2(q) \times \beta_3(r) \times \eta(q \rightarrow r) \times \eta(r \downarrow x_3)}{p(\mathbf{x}) = \alpha_{|\mathbf{x}|}(\text{STOP})}$$



$$\begin{aligned}
 p(y_2 = q, y_3 = r \mid \mathbf{x}) &\propto p(y_2 = q, y_3 = r, \mathbf{x}) \\
 &= \frac{\alpha_2(q) \times \beta_3(r) \times \eta(q \rightarrow r) \times \eta(r \downarrow x_3)}{\sum_{q', r' \in \mathcal{Q}} \alpha_2(q') \times \beta_3(r') \times \eta(q' \rightarrow r') \times \eta(r' \downarrow x_3)} \\
 &= \frac{\alpha_2(q) \times \beta_3(r) \times \eta(q \rightarrow r) \times \eta(r \downarrow x_3)}{p(\mathbf{x}) = \alpha_{|\mathbf{x}|}(\text{STOP})}
 \end{aligned}$$



$$\begin{aligned}
 p(y_2 = q, y_3 = r \mid \mathbf{x}) &\propto p(y_2 = q, y_3 = r, \mathbf{x}) \\
 &= \frac{\alpha_2(q) \times \beta_3(r) \times \eta(q \rightarrow r) \times \eta(r \downarrow x_3)}{\sum_{q', r' \in \mathcal{Q}} \alpha_2(q') \times \beta_3(r') \times \eta(q' \rightarrow r') \times \eta(r' \downarrow x_3)} \\
 &= \frac{\alpha_2(q) \times \beta_3(r) \times \eta(q \rightarrow r) \times \eta(r \downarrow x_3)}{p(\mathbf{x}) = \alpha_{|\mathbf{x}|}(\text{STOP})}
 \end{aligned}$$

$$\begin{aligned}
p(y_2 = q, y_3 = r \mid \mathbf{x}) &\propto p(y_2 = q, y_3 = r, \mathbf{x}) \\
&= \frac{\alpha_2(q) \times \beta_3(r) \times \eta(q \rightarrow r) \times \eta(r \downarrow x_3)}{\sum_{q', r' \in \mathcal{Q}} \alpha_2(q') \times \beta_3(r') \times \eta(q' \rightarrow r') \times \eta(r' \downarrow x_3)} \\
&= \frac{\alpha_2(q) \times \beta_3(r) \times \eta(q \rightarrow r) \times \eta(r \downarrow x_3)}{p(\mathbf{x}) = \alpha_{|\mathbf{x}|}(\text{STOP})}
\end{aligned}$$

The expectation over the full structure is then

$$\mathbb{E}[\text{freq}(q \rightarrow r)] = \sum_{i=1}^{|\mathbf{x}|} p(y_i = q, y_{i+1} = r \mid \mathbf{x})$$

$$\begin{aligned}
p(y_2 = q, y_3 = r \mid \mathbf{x}) &\propto p(y_2 = q, y_3 = r, \mathbf{x}) \\
&= \frac{\alpha_2(q) \times \beta_3(r) \times \eta(q \rightarrow r) \times \eta(r \downarrow x_3)}{\sum_{q', r' \in \mathcal{Q}} \alpha_2(q') \times \beta_3(r') \times \eta(q' \rightarrow r') \times \eta(r' \downarrow x_3)} \\
&= \frac{\alpha_2(q) \times \beta_3(r) \times \eta(q \rightarrow r) \times \eta(r \downarrow x_3)}{p(\mathbf{x}) = \alpha_{|\mathbf{x}|}(\text{STOP})}
\end{aligned}$$

The expectation over the full structure is then

$$\mathbb{E}[\text{freq}(q \rightarrow r)] = \sum_{i=1}^{|\mathbf{x}|} p(y_i = q, y_{i+1} = r \mid \mathbf{x})$$

The expectation over state occupancy is

$$\mathbb{E}[\text{freq}(q)] = \sum_{r \in \mathcal{Q}} \mathbb{E}[\text{freq}(q \rightarrow r)]$$

$$\begin{aligned}
p(y_2 = q, y_3 = r \mid \mathbf{x}) &\propto p(y_2 = q, y_3 = r, \mathbf{x}) \\
&= \frac{\alpha_2(q) \times \beta_3(r) \times \eta(q \rightarrow r) \times \eta(r \downarrow x_3)}{\sum_{q', r' \in \mathcal{Q}} \alpha_2(q') \times \beta_3(r') \times \eta(q' \rightarrow r') \times \eta(r' \downarrow x_3)} \\
&= \frac{\alpha_2(q) \times \beta_3(r) \times \eta(q \rightarrow r) \times \eta(r \downarrow x_3)}{p(\mathbf{x}) = \alpha_{|\mathbf{x}|}(\text{STOP})}
\end{aligned}$$

The expectation over the full structure is then

$$\mathbb{E}[\text{freq}(q \rightarrow r)] = \sum_{i=1}^{|\mathbf{x}|} p(y_i = q, y_{i+1} = r \mid \mathbf{x})$$

The expectation over state occupancy is

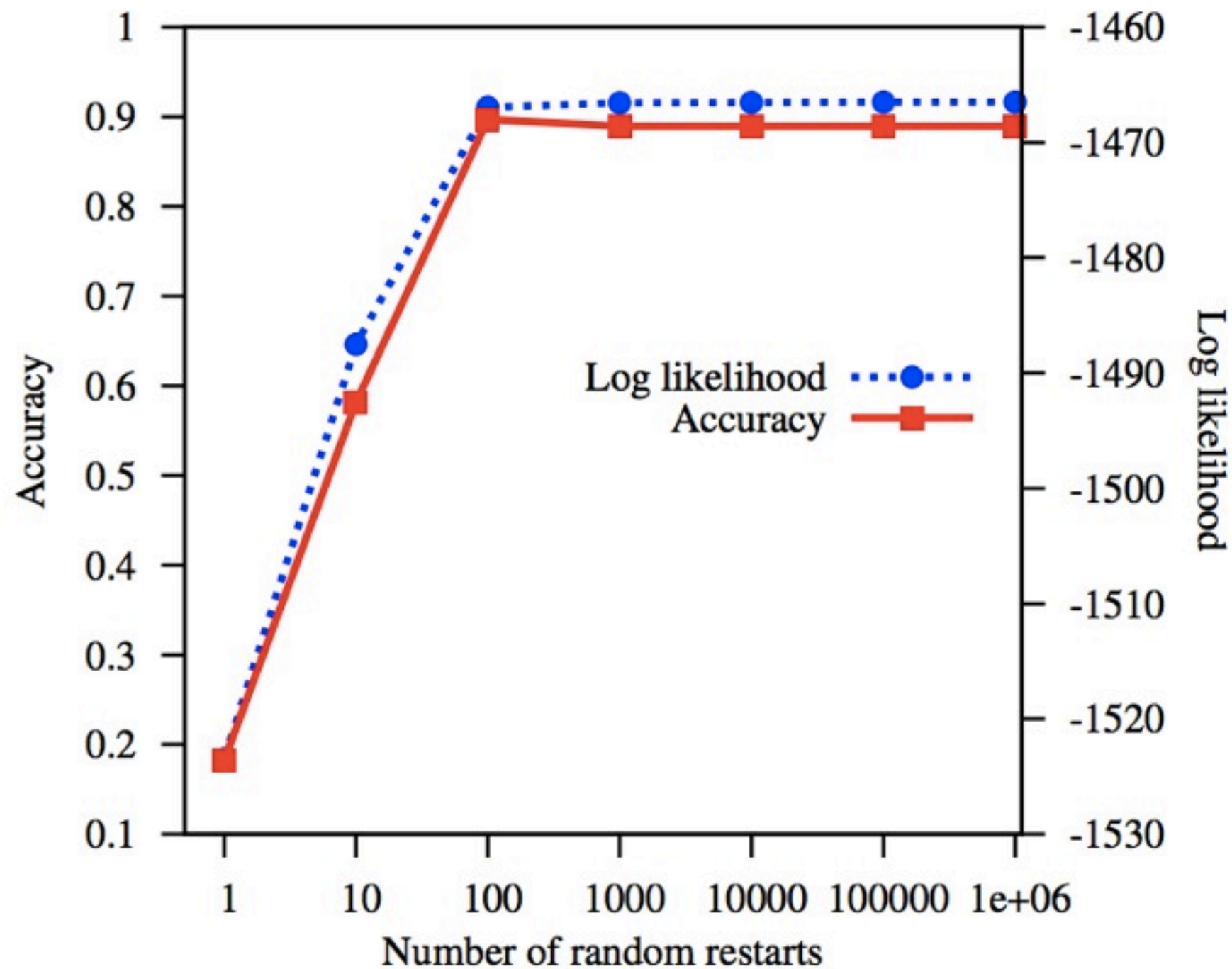
$$\mathbb{E}[\text{freq}(q)] = \sum_{r \in \mathcal{Q}} \mathbb{E}[\text{freq}(q \rightarrow r)]$$

What is $\mathbb{E}[\text{freq}(q \downarrow x)]$?

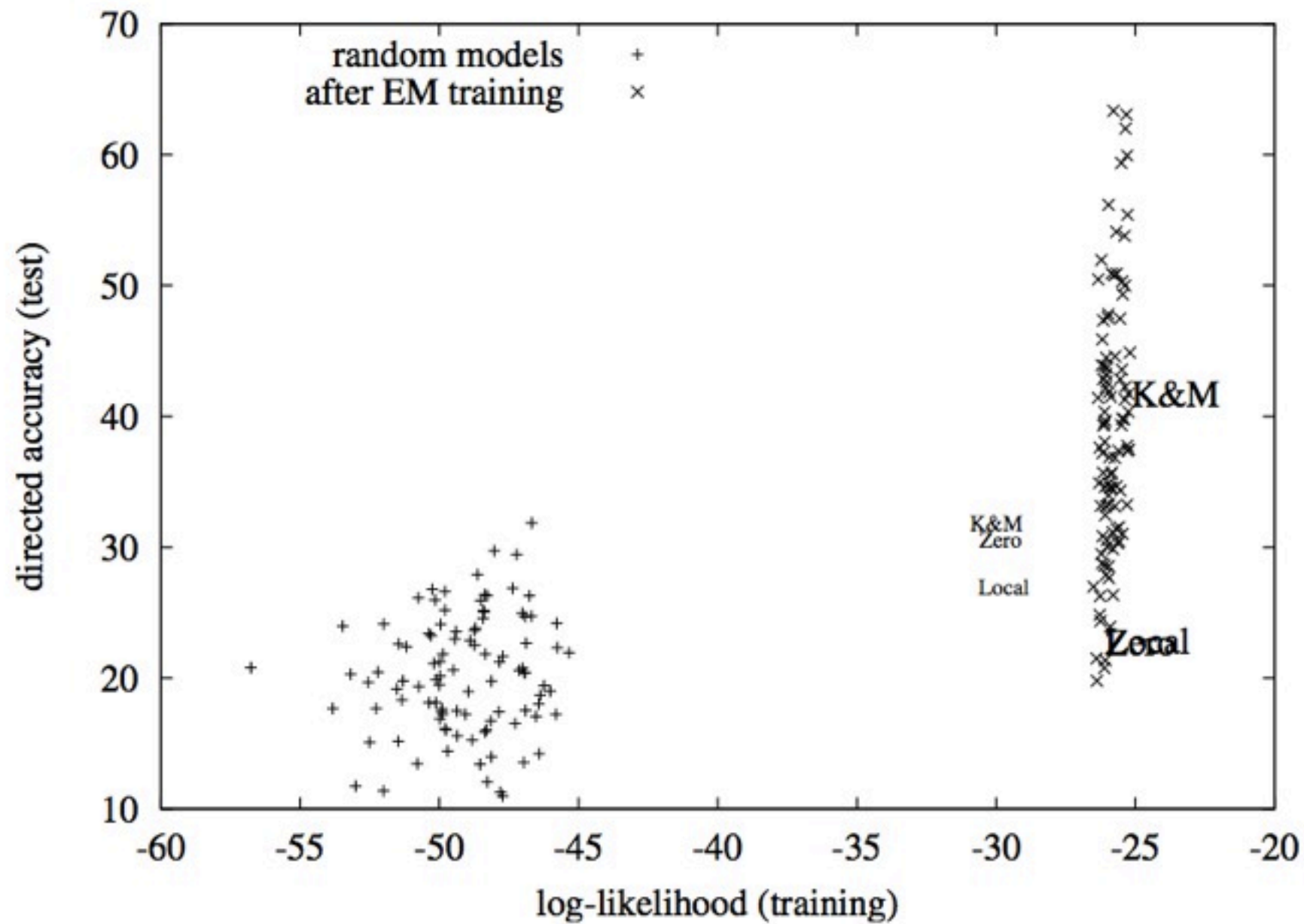
Random Restarts

- Non-convex optimization only finds a local solution
- Several strategies
 - Random restarts
 - Simulated annealing

Decipherment



Grammar Induction



Inductive Bias

- A model can learn nothing without inductive bias ... whence inductive bias?
 - Model structure
 - Priors (next week)
 - Posterior regularization (Google it)
- **Features** provide a very flexible means to bias a model

EM with Features

- Let's replace the multinomials with log linear distributions

$$\begin{aligned}\eta(q \rightarrow r) &= \theta_{q,r} \\ &= \frac{\exp \mathbf{w}^\top \mathbf{f}(q, r)}{\sum_{q' \in \mathcal{Q}} \exp \mathbf{w}^\top \mathbf{f}(q', r)}\end{aligned}$$

How will the likelihood of this model compare to the likelihood of the previous model?

EM with Features

- Let's replace the multinomials with log linear distributions

$$\begin{aligned}\eta(q \rightarrow r) &= \theta_{q,r} \\ &= \frac{\exp \mathbf{w}^\top \mathbf{f}(q, r)}{\sum_{q' \in \mathcal{Q}} \exp \mathbf{w}^\top \mathbf{f}(q', r)}\end{aligned}$$

How will the likelihood of this model compare to the likelihood of the previous model?

Learning Algorithm I

- E step
 - given model parameters, compute posterior distribution over transitions (states, etc)
 - compute $\mathbb{E}_{q(\mathbf{y})} \sum_{q,r} f(q,r)$
 - These are your “empirical” expectations

Learning Algorithm I

- M step
- The gradient of the expected log likelihood of \mathbf{x}, \mathbf{y} under $q(\mathbf{y})$ is

$$\nabla \mathbb{E}_{q(\mathbf{y})} \log p(\mathbf{x}, \mathbf{y}) = \mathbb{E}_q \sum_{q,r} \mathbf{f}(q, r) - \sum_{q,r} \mathbb{E}_q[\text{freq}(q)] \mathbb{E}_{p(r|q;\mathbf{w})} \mathbf{f}(q, r)$$

- Use LBFGS or gradient descent to solve