# Lecture 4.3: Learning lexical translations

## IBM model 1
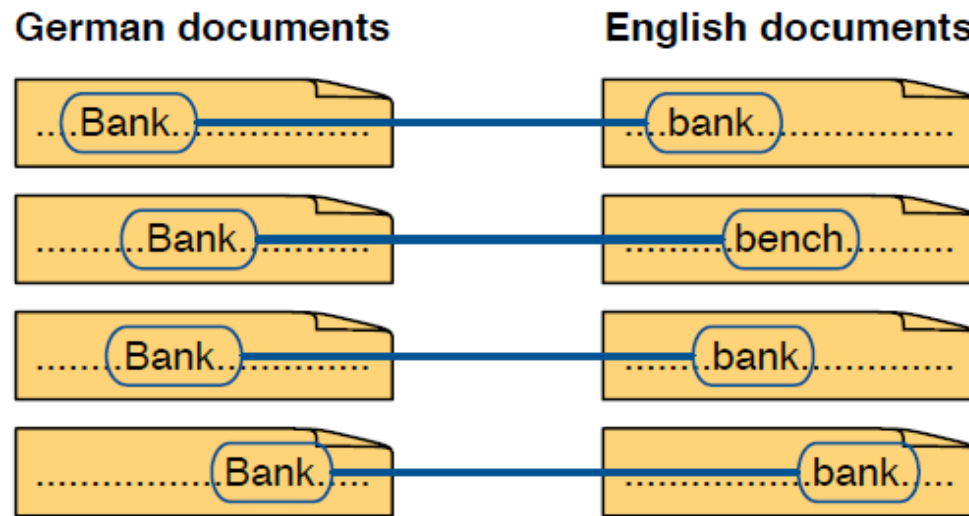
*Preparation for the lab*
*slides borrowed from Philipp Koehn*

# Guessing game:

| Indonesian | English | |
|------------|---------|---|
| Ayam bakar | Grilled chicken |  |
| Ayam goreng | Fried chicken |  |
| Bebek goreng | Fried duck |  |
| Ikan bakar | Grilled fish |  |
| Ikan goreng | ? ? |  |

Goreng=?
Ayam=?
Bakar=?
Ikan=?
Bebek=?

- Collect counts, infer probabilities

**German documents**      **English documents**

...(.Bank.)................. ————— ...bank.................

........(..Bank.)........... ————— .........bench.........

.......(.Bank.)............. ————— .......bank.............

................(.Bank.).... ————— ................bank....
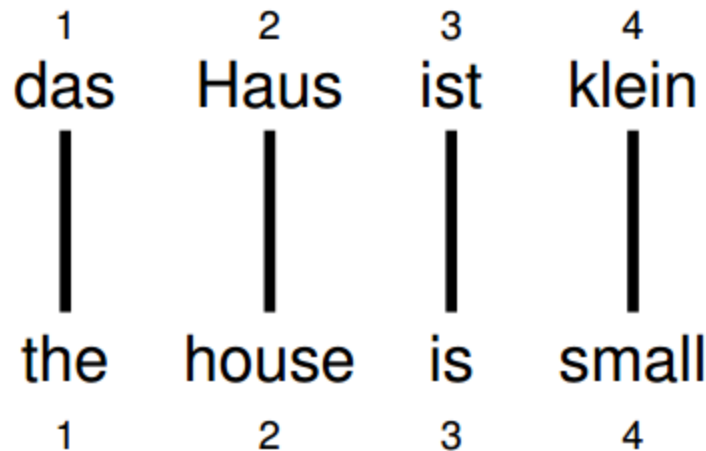
⇒ p(bank|Bank) = 0.75, p(bench|Bank) = 0.25

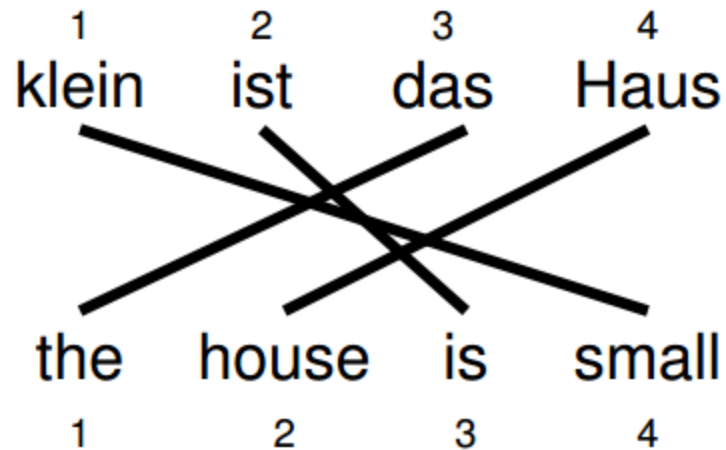Courtesy of Philipp Koehn, Professor at The Johns Hopkins University

# Alignments

- Alignments are the "links" between words in parallel sentences
- Word positions numbered 1-4



English target word at position $i$ to a German source word at position $j$ with a function $a : i \rightarrow j$
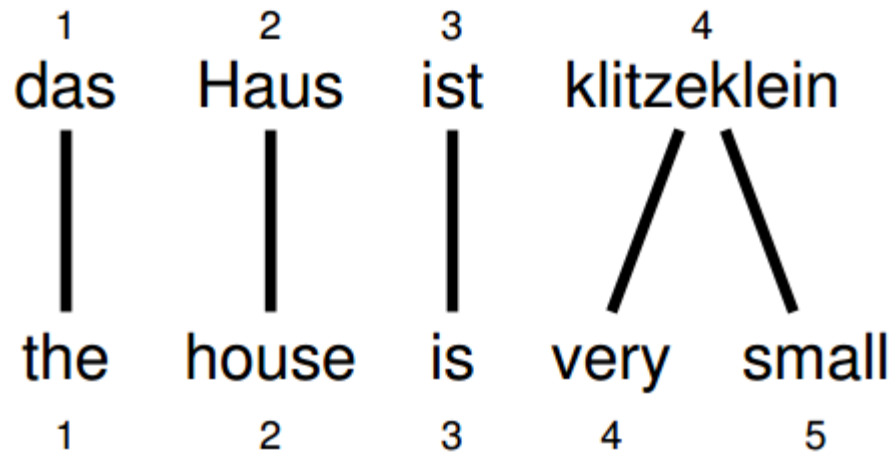
$$a : \{1 \rightarrow 1, 2 \rightarrow 2, 3 \rightarrow 3, 4 \rightarrow 4\}$$
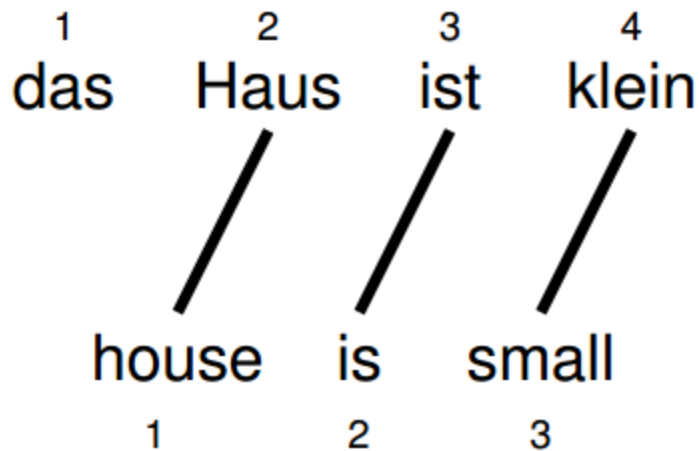
# Alignment reordering



| 1 | 2 | 3 | 4 |
|---|---|---|---|
| klein | ist | das | Haus |

| the | house | is | small |
|---|---|---|---|
| 1 | 2 | 3 | 4 |

$$a : \{1 \rightarrow 3, 2 \rightarrow 4, 3 \rightarrow 2, 4 \rightarrow 1\}$$

# Alignment: one to many

das Haus ist klitzeklein
the house is very small

(1) das (2) Haus (3) ist (4) klitzeklein

(1) the (2) house (3) is (4) very (5) small

$$a : \{1 \to 1, 2 \to 2, 3 \to 3, 4 \to 4, 5 \to 4\}$$

# Dropping/insert words



$a : \{1 \rightarrow 2, 2 \rightarrow 3, 3 \rightarrow 4\}$

$a : \{1 \rightarrow 1, 2 \rightarrow 2, 3 \rightarrow 3, 4 \rightarrow 0, 5 \rightarrow 4\}$

# IBM model 1

- Generative model: break up translation process into smaller steps

  IBM Model 1 only uses lexical translation

- Translation probability
  - for a foreign sentence $f = (f_1, ..., f_{lf})$ of length $lf$
  - to an English sentence $e = (e_1, ..., e_{le})$ of length $le$
  - with an alignment of each English word $e_j$ to a foreign word $f_i$ according to the alignment function $a : j \rightarrow i$
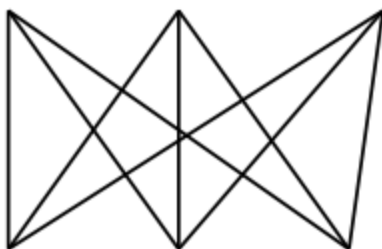
$$p(\mathbf{e}, a | \mathbf{f}) = \frac{\epsilon}{(l_f + 1)^{l_e}} \prod_{j=1}^{l_e} t(e_j | f_{a(j)})$$

- parameter ε is a normalization constant

- Incomplete data
  - if we had complete data, we could estimate model
  - if we had model, we could fill in the gaps in the data

- Expectation Maximization (EM) in a nutshell
  1. initialize model parameters (e.g. uniform)
  2. assign probabilities to the missing data
  3. estimate model parameters from completed data
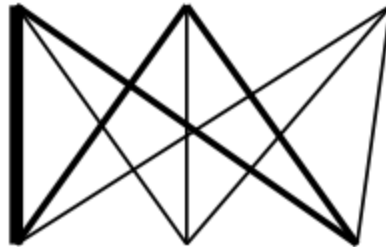  4. iterate steps 2–3 until convergence

... 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

... 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

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



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

- After another iteration

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

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

- Convergence

- Inherent hidden structure revealed by EM

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



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

```
p(la|the) = 0.453
p(le|the) = 0.334
p(maison|house) = 0.876
p(bleu|blue) = 0.563
...
```

- Parameter estimation from the aligned corpus

- **Probabilities**

$$p(\text{the}|\text{la}) = 0.7 \qquad p(\text{house}|\text{la}) = 0.05$$
$$p(\text{the}|\text{maison}) = 0.1 \qquad p(\text{house}|\text{maison}) = 0.8$$

- **Alignments**



$$p(\mathbf{e}, a|\mathbf{f}) = 0.56 \qquad p(\mathbf{e}, a|\mathbf{f}) = 0.035 \qquad p(\mathbf{e}, a|\mathbf{f}) = 0.08 \qquad p(\mathbf{e}, a|\mathbf{f}) = 0.005$$

$$p(a|\mathbf{e}, \mathbf{f}) = 0.824 \qquad p(a|\mathbf{e}, \mathbf{f}) = 0.052 \qquad p(a|\mathbf{e}, \mathbf{f}) = 0.118 \qquad p(a|\mathbf{e}, \mathbf{f}) = 0.007$$

- **Counts**

$$c(\text{the}|\text{la}) = 0.824 + 0.052 \qquad c(\text{house}|\text{la}) = 0.052 + 0.007$$
$$c(\text{the}|\text{maison}) = 0.118 + 0.007 \qquad c(\text{house}|\text{maison}) = 0.824 + 0.118$$

# IBM model 1 pseudo code

**Input:** set of sentence pairs $(\mathbf{e}, \mathbf{f})$
**Output:** translation prob. $t(e|f)$

1:   initialize $t(e|f)$ uniformly
2:   **while** not converged **do**
3:     // initialize
4:     $\text{count}(e|f) = 0$ **for all** $e, f$
5:     $\text{total}(f) = 0$ **for all** $f$
6:     **for all** sentence pairs $(\mathbf{e},\mathbf{f})$ **do**
7:      // compute normalization
8:      **for all** words $e$ in $\mathbf{e}$ **do**
9:       $\text{s-total}(e) = 0$
10:      **for all** words $f$ in $\mathbf{f}$ **do**
11:        $\text{s-total}(e) \mathrel{+}= t(e|f)$
12:      **end for**
13:    **end for**

14:     // collect counts
15:     **for all** words $e$ in $\mathbf{e}$ **do**
16:      **for all** words $f$ in $\mathbf{f}$ **do**
17:       $\text{count}(e|f) \mathrel{+}= \frac{t(e|f)}{\text{s-total}(e)}$
18:       $\text{total}(f) \mathrel{+}= \frac{t(e|f)}{\text{s-total}(e)}$
19:      **end for**
20:     **end for**
21:   **end for**
22:   // estimate probabilities
23:   **for all** foreign words $f$ **do**
24:     **for all** English words $e$ **do**
25:      $t(e|f) = \frac{\text{count}(e|f)}{\text{total}(f)}$
26:     **end for**
27:   **end for**
28: **end while**

# Convergence

das    Haus          das    Buch          ein    Buch

the    house         the    book          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 |