

# Machine Translation 2: Phrase-Based Translation

Computational Linguistics

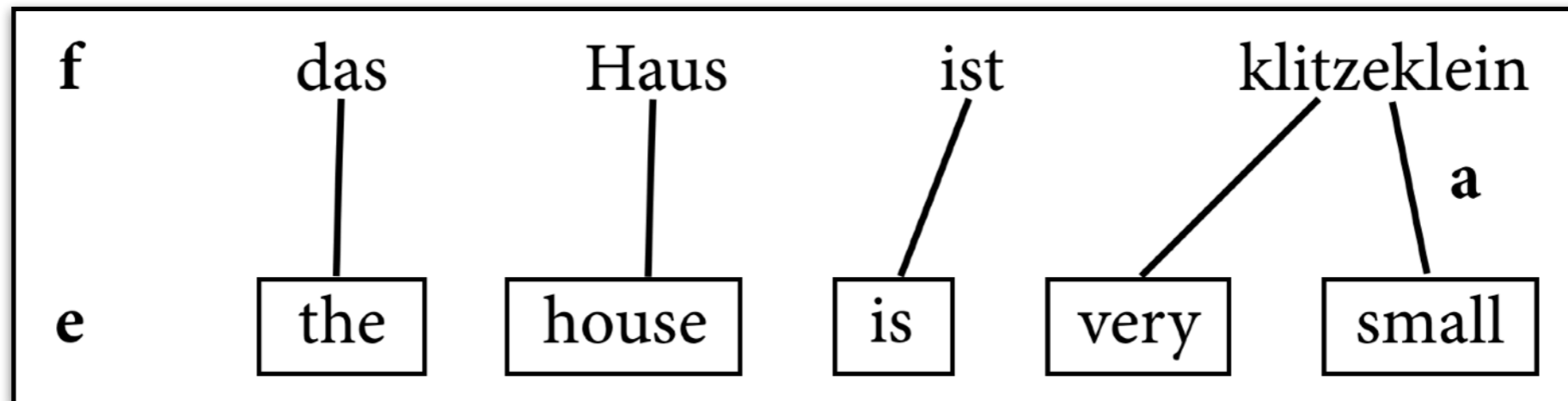
Alexander Koller

7 January 2020

slides contain material from [mt-class.org](http://mt-class.org)

# Where were we?

- Last time: Word alignments.



- Today: Actual machine translation.
  - ▶ Input: “Das Haus ist klitzeklein.”
  - ▶ Output: “The house is very small.”

# Translation quality

- We can measure quality of a translation in two dimensions:
  - ▶ *Adequacy*: How accurately does translation represent the meaning of the original?
  - ▶ *Fluency*: Is the translation a good string of the target language (“good English”)?
- How can we select a fluent translation?

# Fluency

Israeli officials are responsible for airport security.

Israel is in charge of the security at this airport.

The security work for this airport is the responsibility of the Israel government.

Israeli side was in charge of the security of this airport.

Israel is responsible for the airport's security.

Israel is responsible for safety work at this airport.

Israel presides over the security of the airport.

Israel took charge of the airport security.

The safety of this airport is taken charge of by Israel.

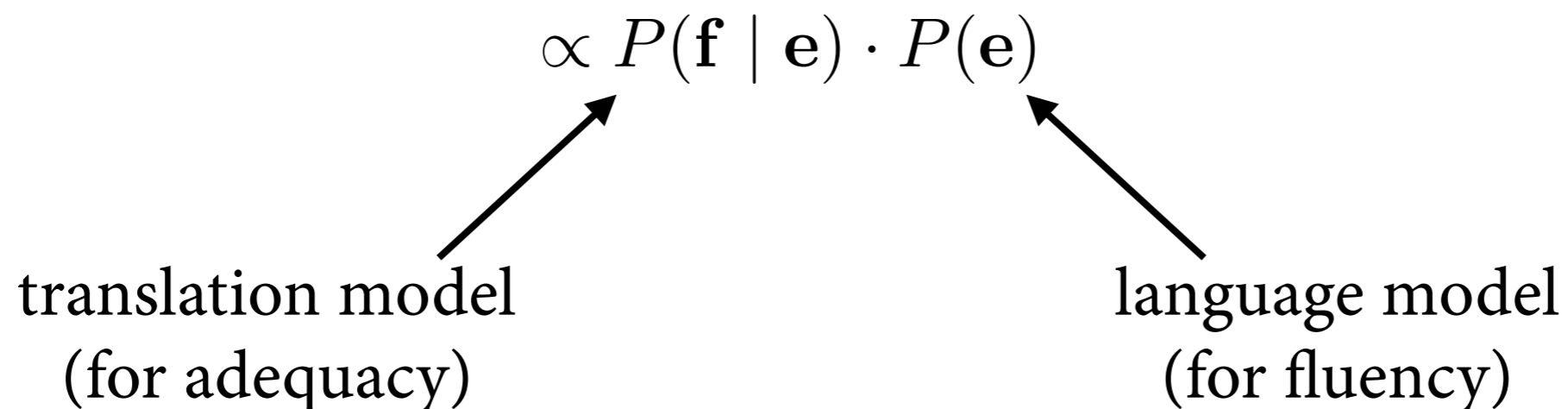
This airport's security is the responsibility of the Israeli security officials.

(from Koehn book)

# Noisy Channel Model

- We can model fluency with a *language model*  $P(\mathbf{e})$  of the target language.
  - ▶ Can estimate from lots of monolingual data!
  - ▶ Use e.g. n-gram models (with smoothing).
- Noisy Channel Model:

$$P(\mathbf{e} \mid \mathbf{f}) = \frac{P(\mathbf{f} \mid \mathbf{e}) \cdot P(\mathbf{e})}{P(\mathbf{f})}$$



# Word-based translation model

- Could derive model for word-by-word translation, e.g. from IBM Model 1:

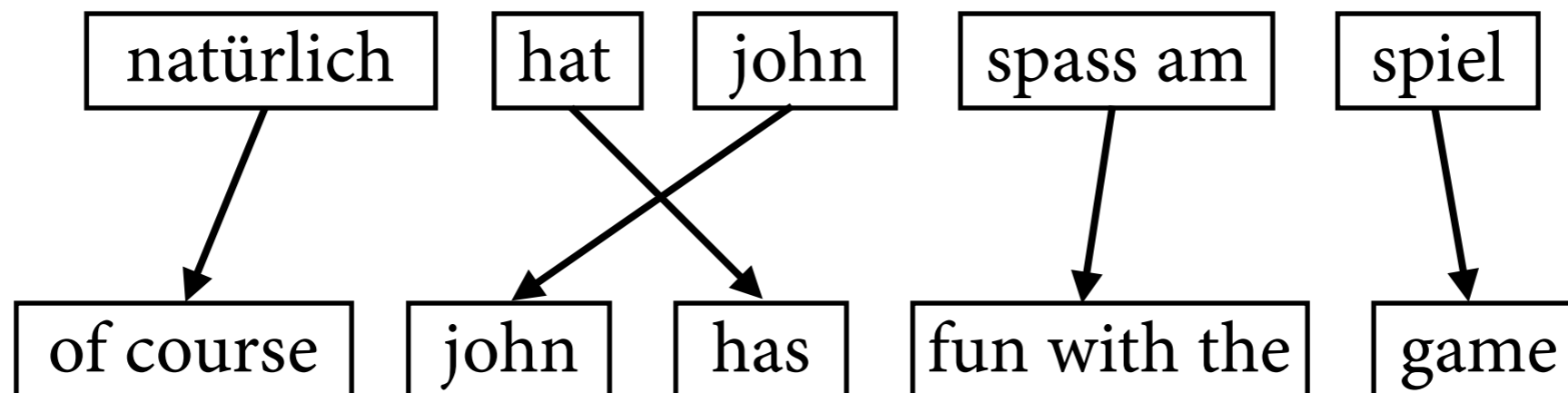
$$P(\mathbf{f} \mid \mathbf{e}) = \sum_{\mathbf{a}} P(\mathbf{f}, \mathbf{a} \mid \mathbf{e})$$

$$\propto \prod_{j=1}^{l_f} \sum_{i=1}^{l_e} P(f_j \mid e_i)$$

- (This would be a terrible translation model.)

# Phrase-based translation

- But want to translate entire *phrases* (i.e. substrings):
  - ▶ translation of one word can consist of multiple words
  - ▶ context of word in phrase can help disambiguate



- Note: these “phrases” need not be linguistically meaningful constituents.

# Phrase-based translation model

The diagram shows the equation for the phrase-based translation model with three annotations and arrows:

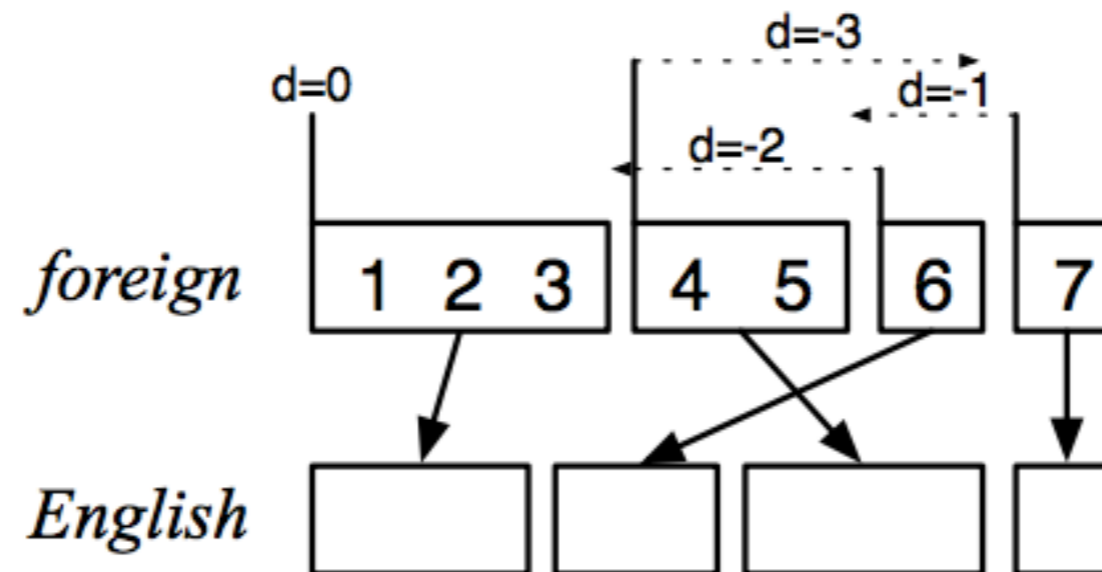
- An arrow from "number of phrases" points to the index  $I$  above the product symbol.
- An arrow from "phrase translation probability" points to the term  $\phi(\bar{f}_i \mid \bar{e}_i)$ .
- An arrow from "distance-based reordering model" points to the term  $d(\text{start}_i - \text{end}_{i-1} - 1)$ .

$$P(\mathbf{f} \mid \mathbf{e}) = \prod_{i=1}^I \phi(\bar{f}_i \mid \bar{e}_i) \cdot d(\text{start}_i - \text{end}_{i-1} - 1)$$

(the whole thing gets multiplied by  $P(\mathbf{e})$  later)

# Reordering Model

Let's assume a simple model for reordering for now.








phrase	translates	movement	distance
1	1–3	start at beginning	0
2	6	skip over 4–5	+2
3	4–5	move back over 4–6	-3
4	7	skip over 6	+1

Scoring function:  $d(x) = \alpha^{|x|}$  — exponential with distance






# Learning phrase translations

- Extend word alignments to phrase alignments.
- Collect all phrase pairs from the parallel corpus (both big and small — we want *all* phrase pairs).
- Estimate phrase translation probabilities  $P(\mathbf{f} \mid \mathbf{e})$  using maximum likelihood estimation (plus smoothing).

# Phrase Extraction






	I	open	the	box
watashi				
wa				
hako				
wo				
akemasu				

# Phrase Extraction

	I	open	the	box
watashi				
wa				
hako				
wo				
akemasu				






akemasu / open

# Phrase Extraction

	I	open	the	box
watashi				
wa				
hako				
wo				
akemasu				






watashi wa / I

# Phrase Extraction

	I	open	the	box
watashi				
wa				
hako				
wo				
akemasu				






hako wo / box

# Phrase Extraction

	I	open	the	box
watashi				
wa				
hako				
wo				
akemasu				

hako wo / the box

# Phrase Extraction

	I	open	the	box
watashi				
wa				
hako				
wo				
akemasu				

hako wo akemasu / open the box

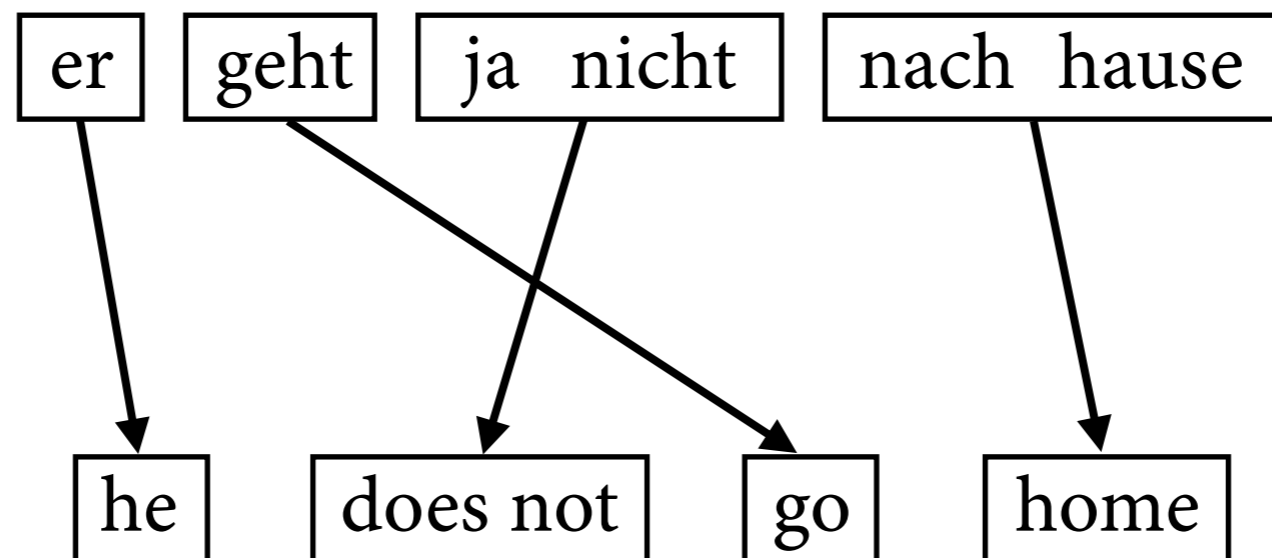
# Decoding

- We now have:
  - ▶ noisy channel  $P(\mathbf{e} \mid \mathbf{f}) \propto P(\mathbf{f} \mid \mathbf{e}) * P(\mathbf{e})$
  - ▶ language model  $P(\mathbf{e})$
  - ▶ phrase-based translation model

$$P(\mathbf{f} \mid \mathbf{e}) = \prod_{i=1}^I \phi(\bar{f}_i \mid \bar{e}_i) \cdot d(\text{start}_i - \text{end}_{i-1} - 1)$$

- We need to solve the *decoding* problem:  
for a given  $\mathbf{f}$ , compute  $\text{argmax}_{\mathbf{e}} P(\mathbf{e} \mid \mathbf{f})$ .

# Basic idea



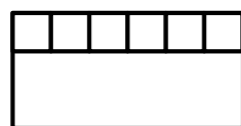
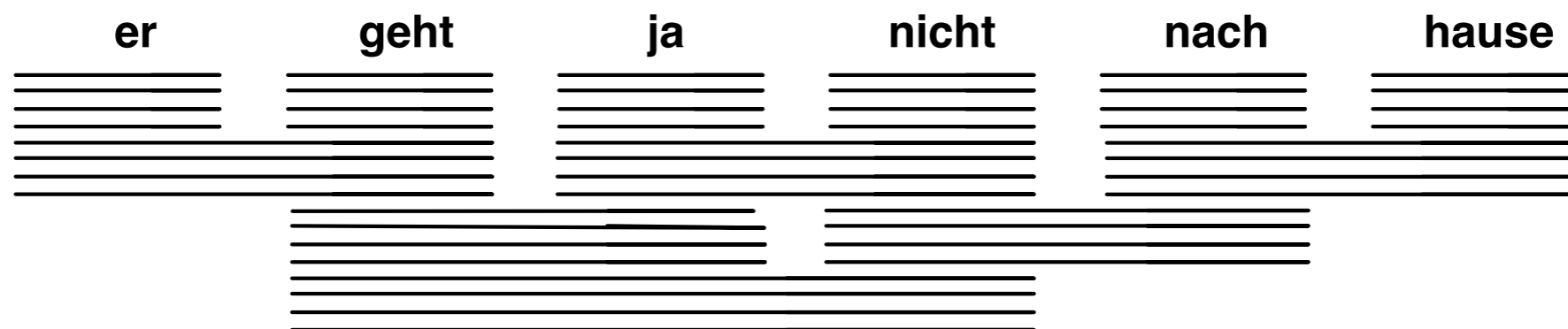
# More realistically

er	geht	ja	nicht	nach	hause
he	is	yes	not	after	house
it	are	is	do not	to	home
, it	goes	, of course	does not	according to	chamber
, he	go	,	is not	in	at home
it is		not		home	
he will be		is not		under house	
it goes		does not		return home	
he goes		do not		do not	
	is		to		
	are		following		
	is after all		not after		
	does		not to		
	not				
	is not				
	are not				
	is not a				

- Many translation options to choose from
  - in Europarl phrase table: 2727 matching phrase pairs for this sentence
  - by pruning to the top 20 per phrase, 202 translation options remain

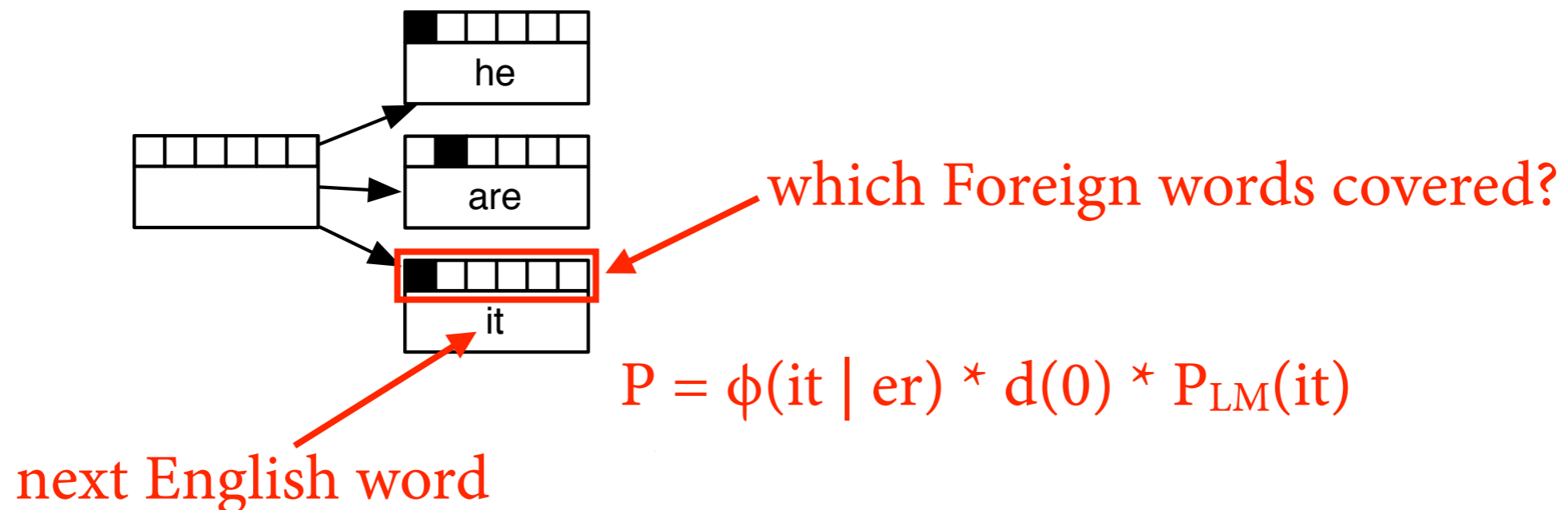
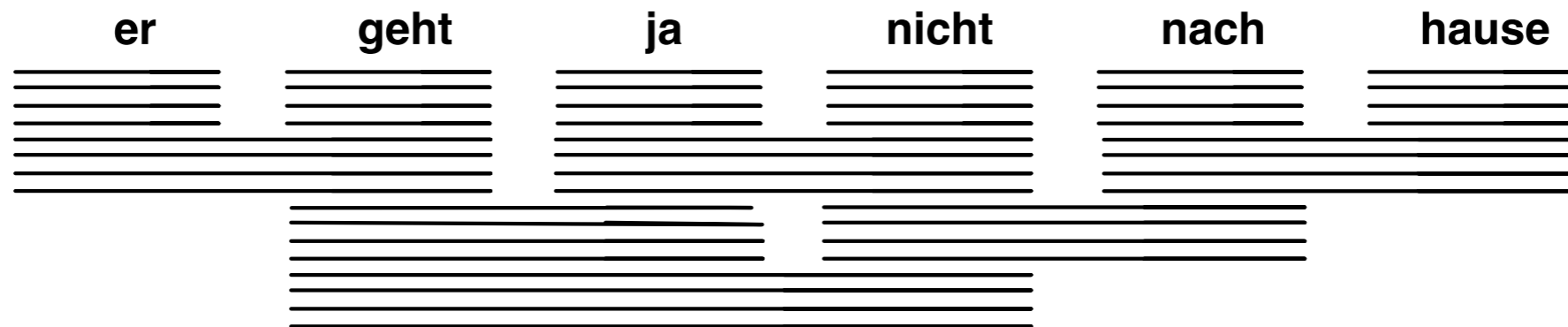
# Decoding as Search

start with empty hypothesis (no words translated)



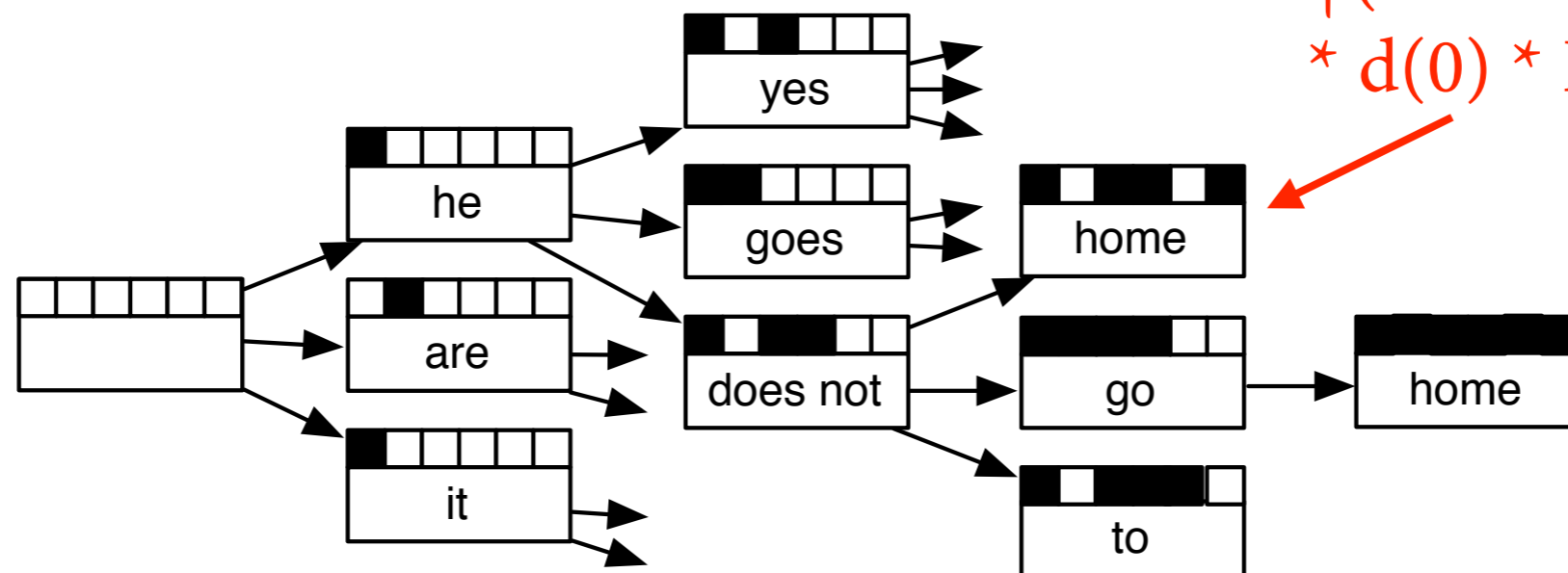
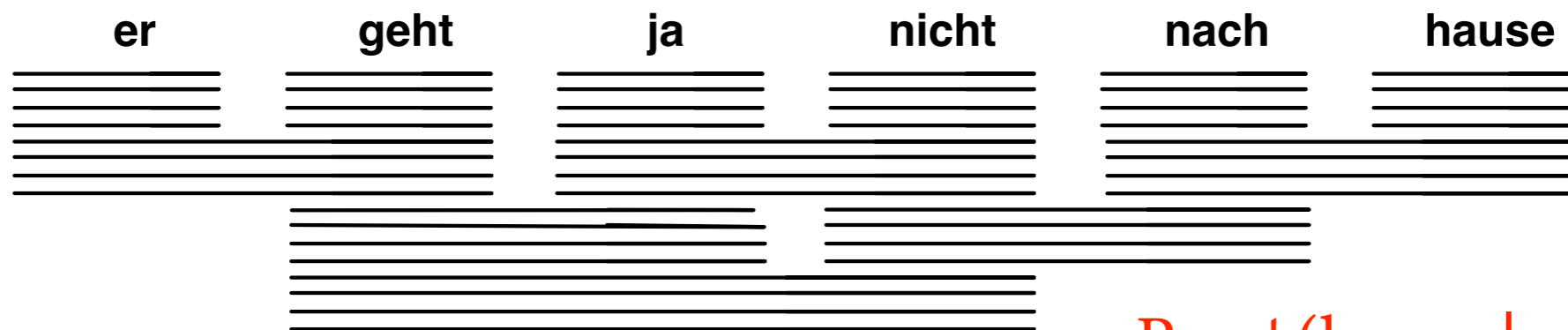
# Decoding as Search

expand hypotheses by next English word



# Decoding as Search

continue expanding hypotheses

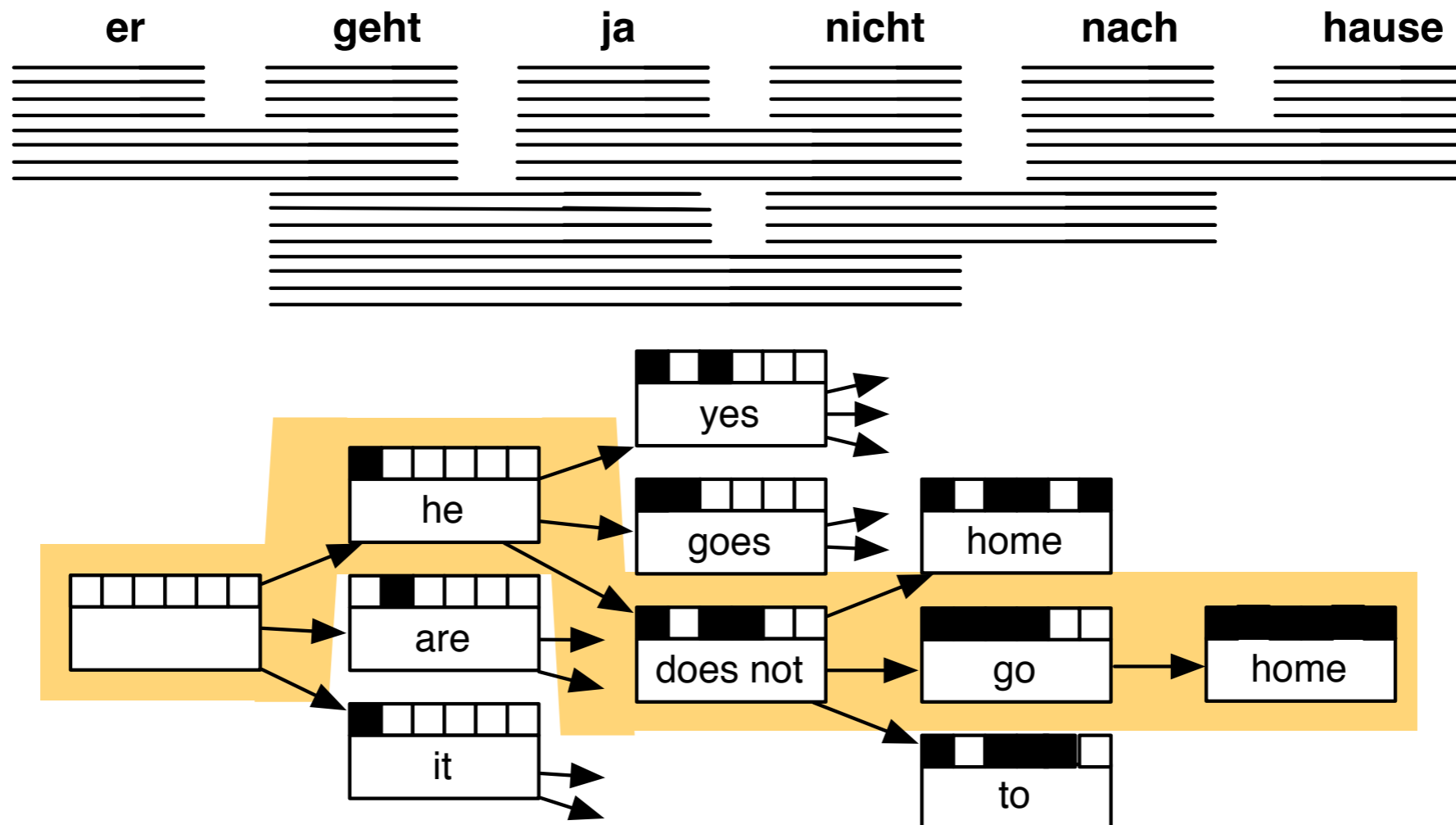


$$P = \phi(\text{home} \mid \text{nach Hause}) \\ * d(0) * P_{\text{LM}}(\text{home} \mid \text{not})$$



# Decoding as Search

backtrack from highest-scoring complete hypothesis



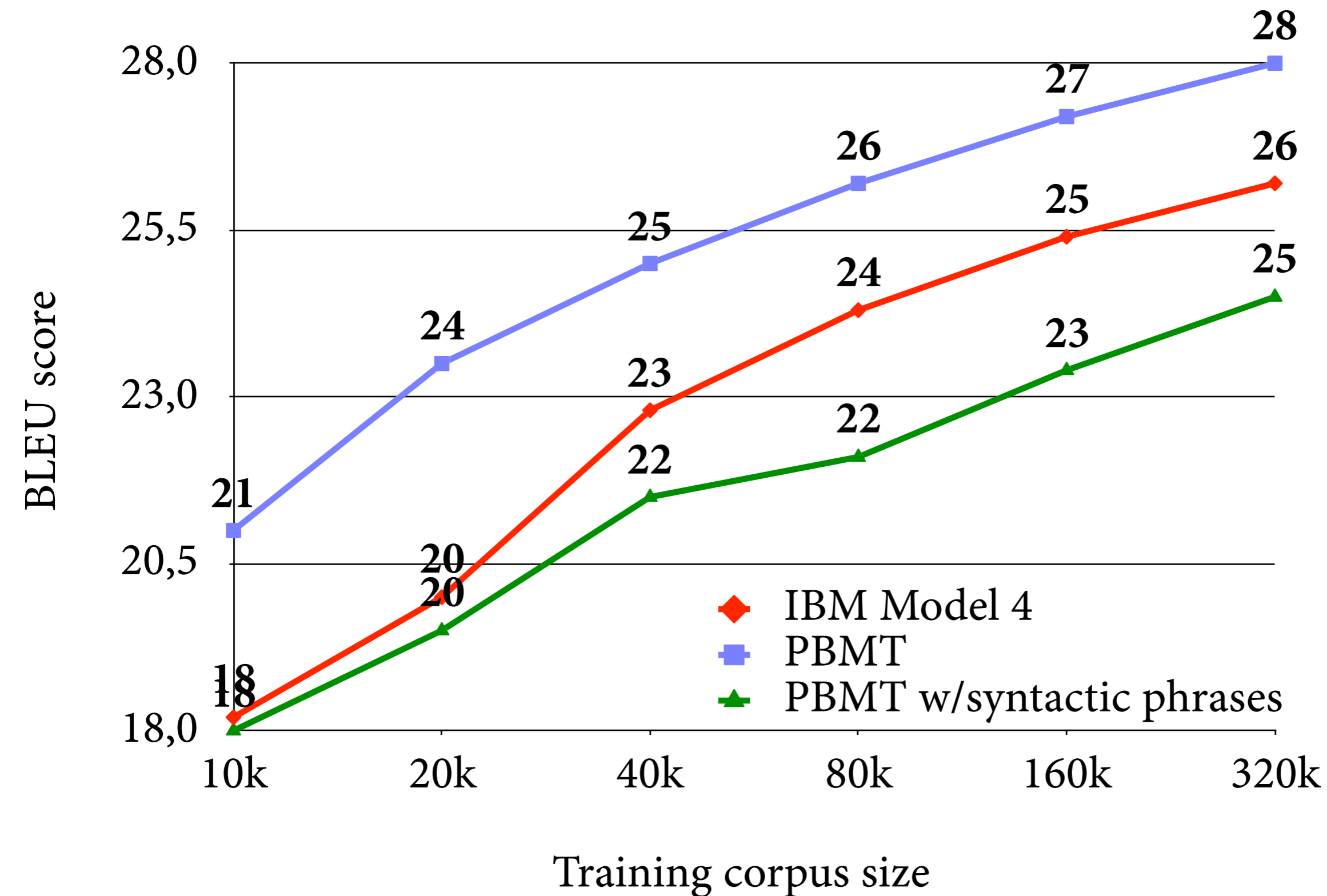
# Computational issues

- Search space is huge.
  - ▶ exponential in sentence length (because of free reordering)
  - ▶ in fact, finding best translation is NP-complete
- Need heuristics to deal with complexity.
  - ▶ beam search: *stack decoding*
  - ▶ A\* search

# Putting linguistics in SMT

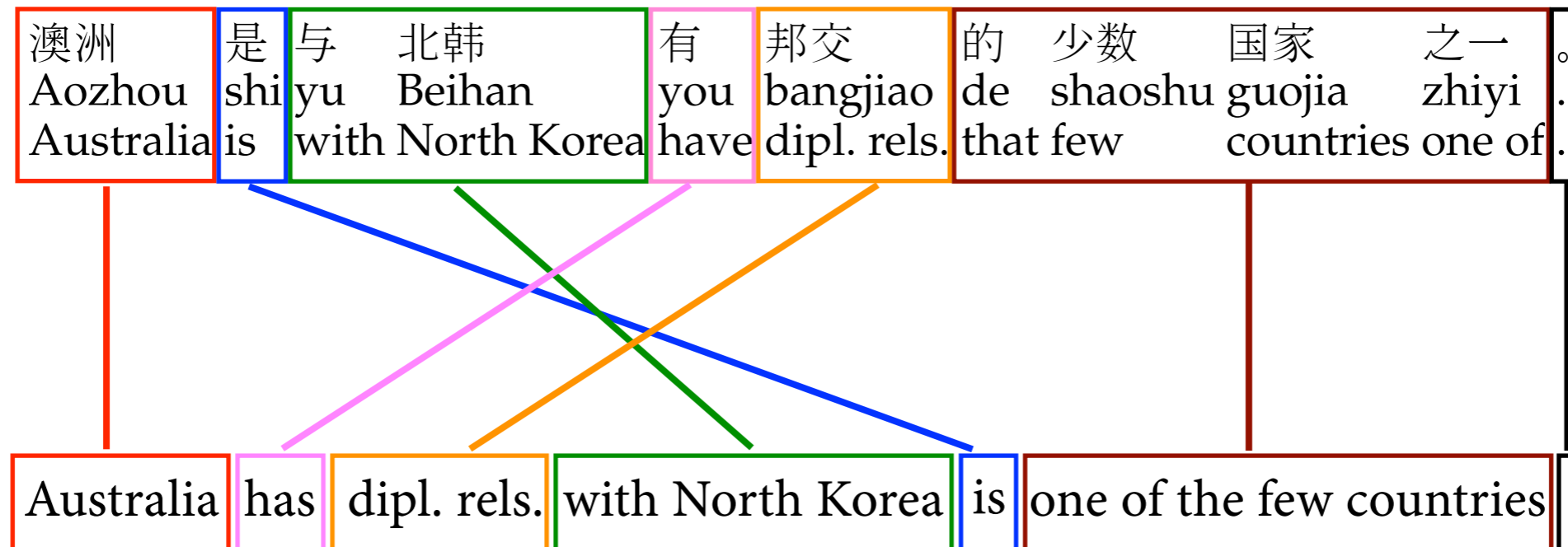
- Word-based, phrase-based SMT very naive from a linguistics perspective.
- Can we do better by putting linguistics into SMT? (At least a bit of syntax?)
- Received wisdom before 2005: phrase-based translation with lots of data much better; syntax hurts.

# Syntax can hurt



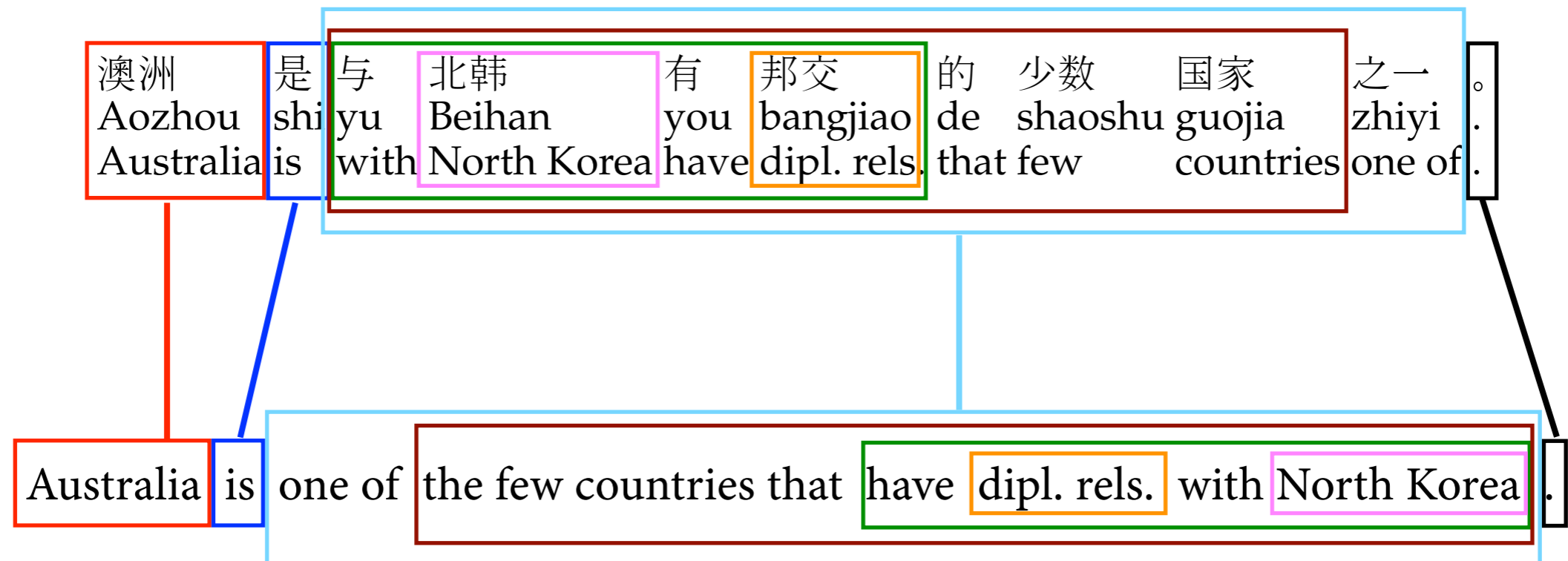
# Chinese-English reordering

(output of phrase-based system ATS)



“Australia is one of the few countries that have diplomatic relations with North Korea.”

# Syntax-based reordering



⟨yu [1] you [2], have [2] with [1]⟩

⟨[1] de [2], the [2] that [1]⟩

⟨[1] zhiyi, one of [1]⟩

“Australia is one of the few countries that have diplomatic relations with North Korea.”

# Syntax-based translation



- Idea: Learn *synchronous* syntax rules that capture syntactic reordering between the two languages.
- Then much less unsystematic reordering necessary.
- We need to figure out:
  - ▶ how to represent translation rules
  - ▶ how to extract translation rules from data
  - ▶ how to define probability model (skipped here)
  - ▶ how to do decoding

# Synchronous CFG

$S \rightarrow X① / X①$

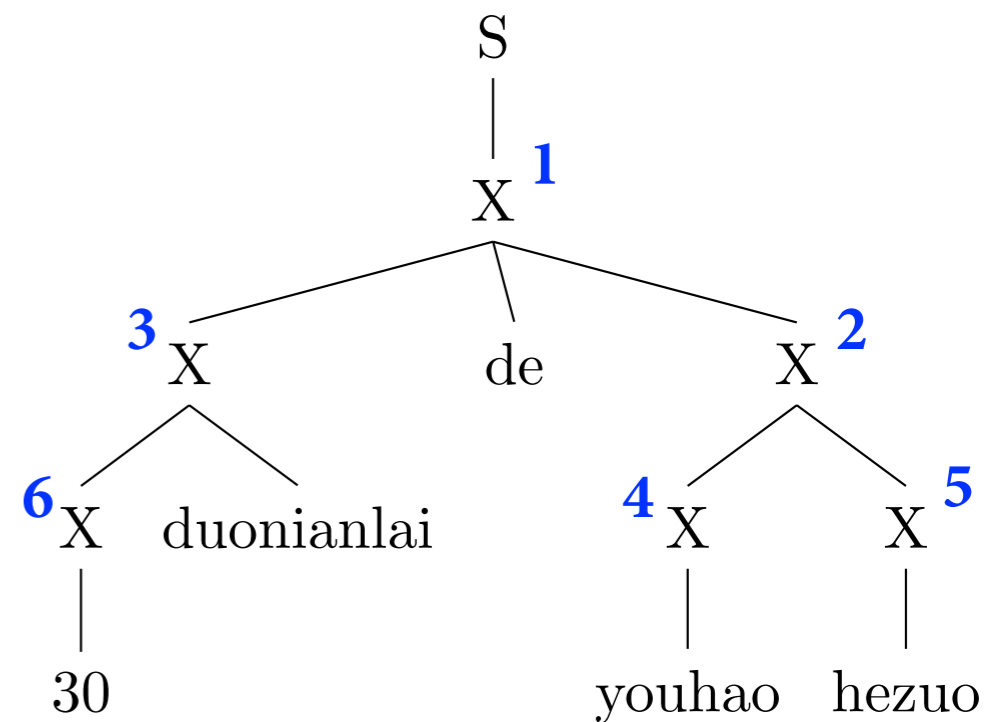
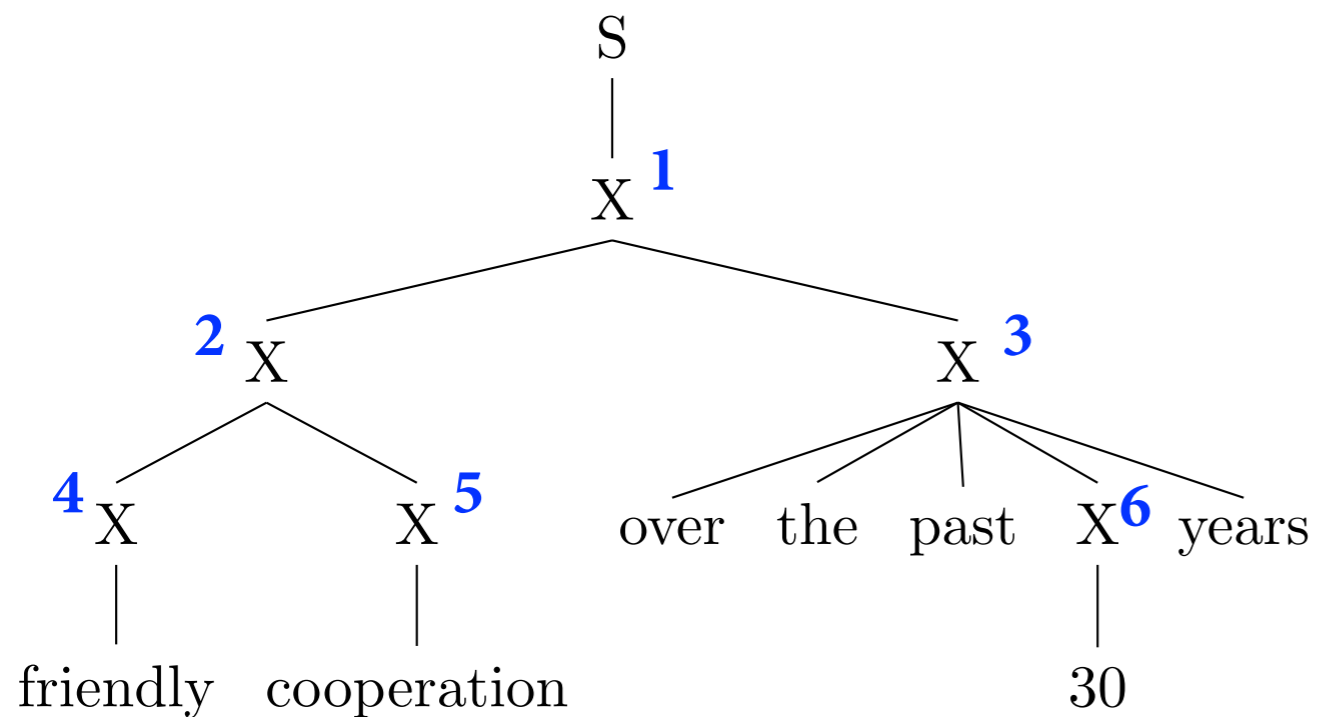
$X \rightarrow X① \text{ de } X② / X② X①$

$X \rightarrow X① X② / X① X②$

$X \rightarrow X① \text{ duonianlai} / \text{over the last } X① \text{ years}$

$X \rightarrow \text{yuohao} / \text{friendly}$

$X \rightarrow 30 / 30$



# SCFG rule extraction

	friendly	cooperation	over	the	last	30	years
30							
duonianlai							
de							
youhao							
hezuo							

$X \rightarrow \text{yuohao} / \text{friendly}$

$X \rightarrow 30 \text{ duonianlai} / \text{over the last 30 years}$

$X \rightarrow 30 / 30$

$X \rightarrow X^{(1)} \text{ duonianlai} / \text{over the last } X^{(1)} \text{ years}$

$X \rightarrow X^{(1)} X^{(2)} / X^{(1)} X^{(2)}$

$X \rightarrow X^{(1)} \text{ de } X^{(2)} / X^{(2)} X^{(1)}$

- Extract all phrase pairs as usual.
- Generate more rules by replacing sub-phrases by nonterminal X.
- Add “glue rules”  $S \rightarrow S^{(1)} X^{(2)} / S^{(1)} X^{(2)}$  and  $S \rightarrow X^{(1)} / X^{(1)}$  to start derivations.

# Decoding schema

$f = \text{"30 duonianlai de youhao hezuo"}$

$[X, 0, 2]$                        $[X, 3, 5]$                        $X \rightarrow X① \text{ de } X② / X② X①$   
*over ... years*                      *friendly ... coop.*

---

$[X, 0, 5]$   
*friendly ... years*

$$\text{prob} = p_1 * p_2 * P(\text{rule}) * P_{\text{LM}}(\text{over} \mid \text{coop.})$$

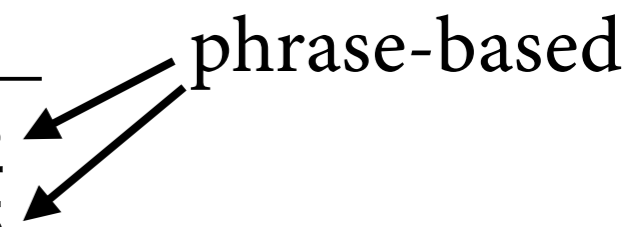
	friendly	cooperation	over	the	last	30	years
30							
duonianlai							
de							
youhao							
hezuo							

# Pruning

- Problem: number of items blown up by factor of  $|V|^{2m}$  for an  $m$ -gram language model.
- Tackle with beam search: for each  $[X, i, j]$  for Foreign positions  $i, j$ , keep only the  $k$  best analyses.
- *Cube pruning*: improve runtime further by only computing a subset of the top  $k$  best analyses for each item by need.

# BLEU Comparison

System	MT03	MT04	MT05
Hiero Monotone	28.27 $\pm$ 1.03	28.83 $\pm$ 0.74	26.35 $\pm$ 0.92
ATS	30.84 $\pm$ 0.99	31.74 $\pm$ 0.73	30.50 $\pm$ 0.95
Hiero	33.72 $\pm$ 1.12	34.57 $\pm$ 0.82	31.79 $\pm$ 0.91



phrase-based

# Results: BLEU and Speed

Method	Settings	Time	BLEU
rescore	$k = 10^4$	16	33.31
rescore	$k = 10^5$	139	33.33
intersect*		1455	37.09
cube prune	$\varepsilon = 0$	23	36.14
cube prune	$\varepsilon = 0.1$	35	36.77
cube prune	$\varepsilon = 0.2$	111	36.91

time in seconds per sentence

# Conclusion

- Noisy channel translation: combine translation model with language model.
- Phrase-based translation: Extract phrases (= arbitrary substrings) from word alignments.
  - ▶ different reordering models, e.g. with SCFGs
- Decoding algorithms must deal with huge search space. Need to do some clever form of beam search.
- Much current research uses neural networks instead.