

What's New in Statistical Machine Translation

Kevin Knight and Philipp Koehn

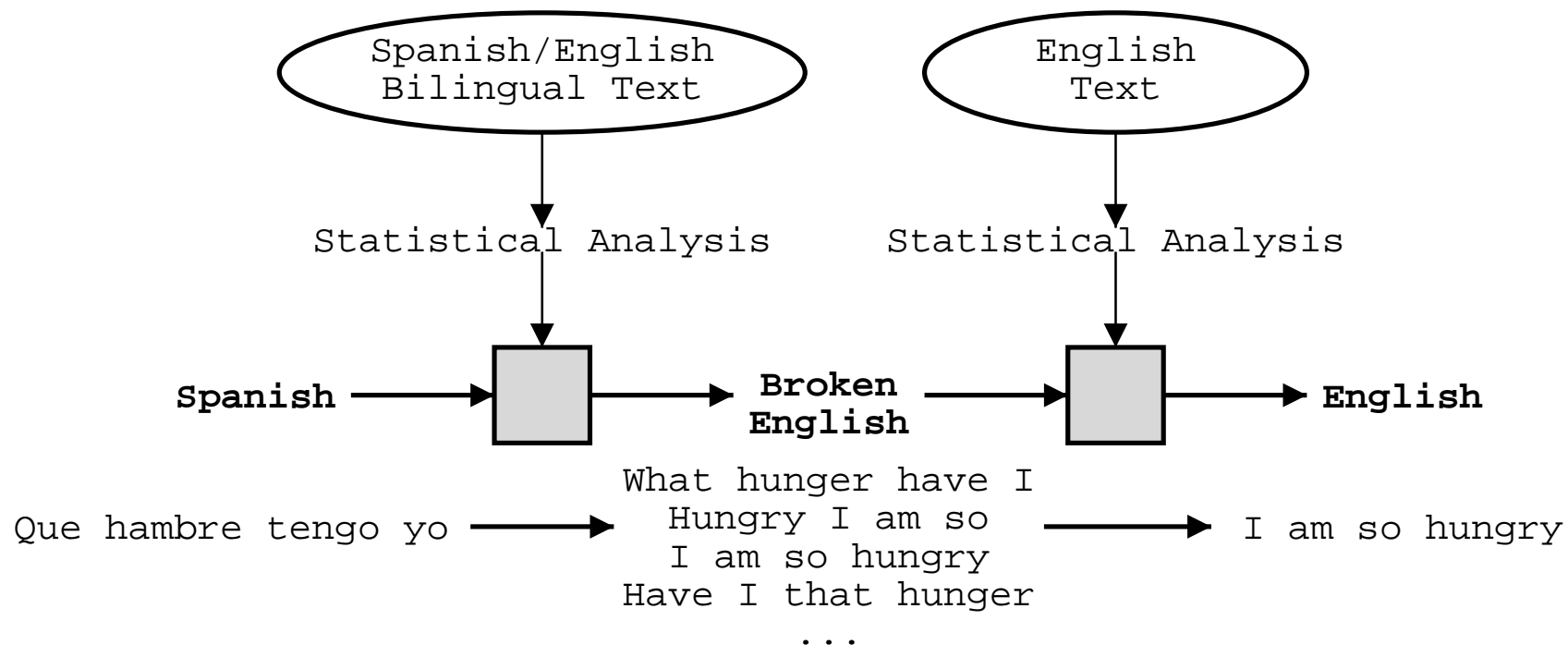
`knight@isi.edu koehn@isi.edu`

Information Sciences Institute
University of Southern California

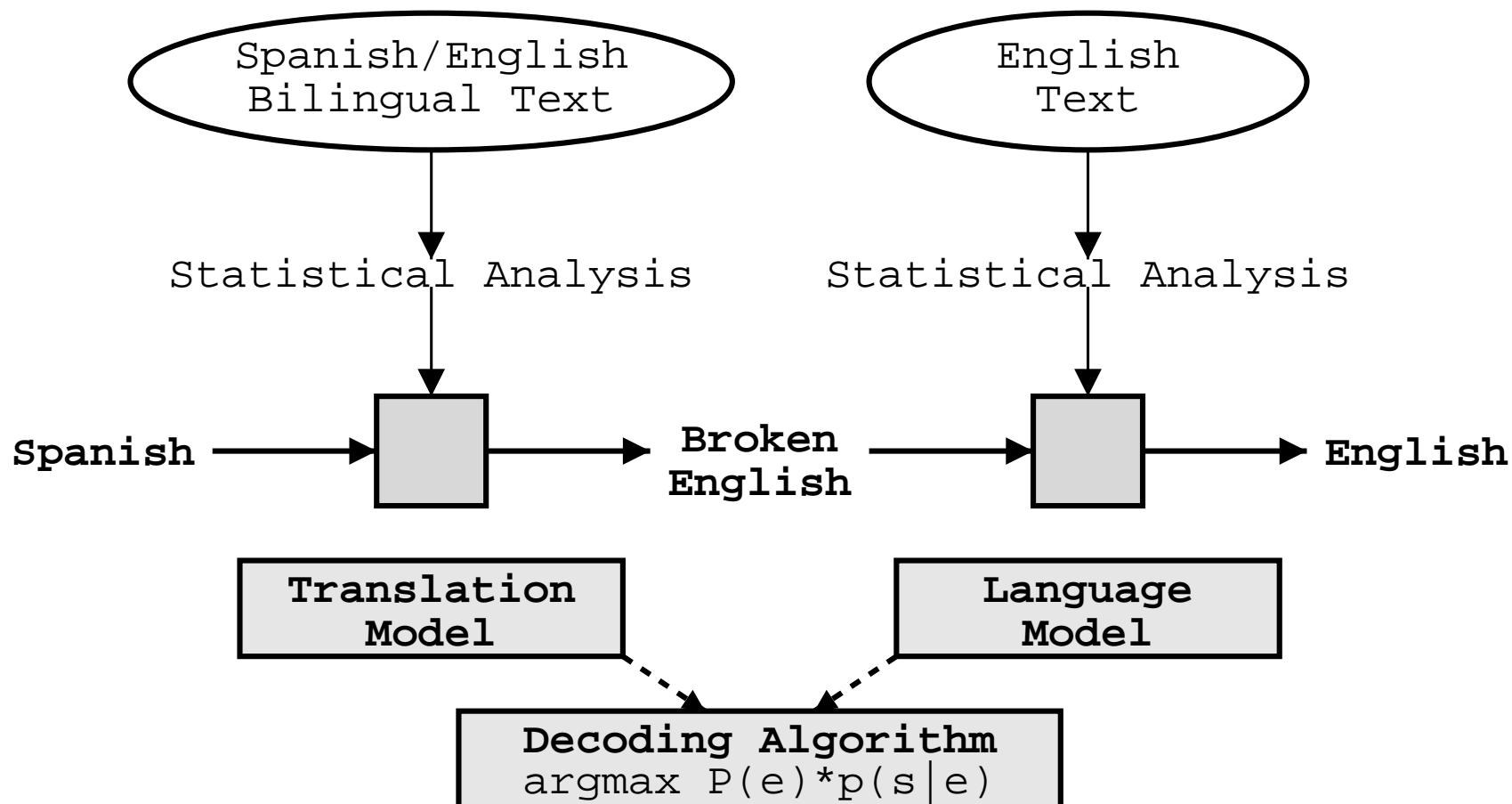
Outline

- Data
- Evaluation
- **Introduction to Statistical Machine Translation**
- Translation Model
- Language Model
- Decoding Algorithm
- New Directions: Divide and Conquer
- Available Resources

Statistical MT Systems



Statistical MT Systems (2)



Three Problems in Statistical MT

- Language Model

- given an English string e , assigns $P(e)$ by formula
- good English string \Rightarrow high $P(e)$
- bad English string \Rightarrow low $P(e)$

- Translation Model

- given a pair of strings $\langle f, e \rangle$, assigns $P(f|e)$ by formula
- $\langle f, e \rangle$ look like translations \Rightarrow high $P(f|e)$
- $\langle f, e \rangle$ don't look like translations \Rightarrow low $P(f|e)$

- Decoding Algorithm

- given a language model, a translation model and a new sentence f ,
find translation e maximizing $P(e) \times P(f|e)$

Outline

- Data
- Evaluation
- Introduction to Statistical Machine Translation
- **Translation Model**
- Language Model
- Decoding Algorithm
- New Directions: Divide and Conquer
- Available Resources

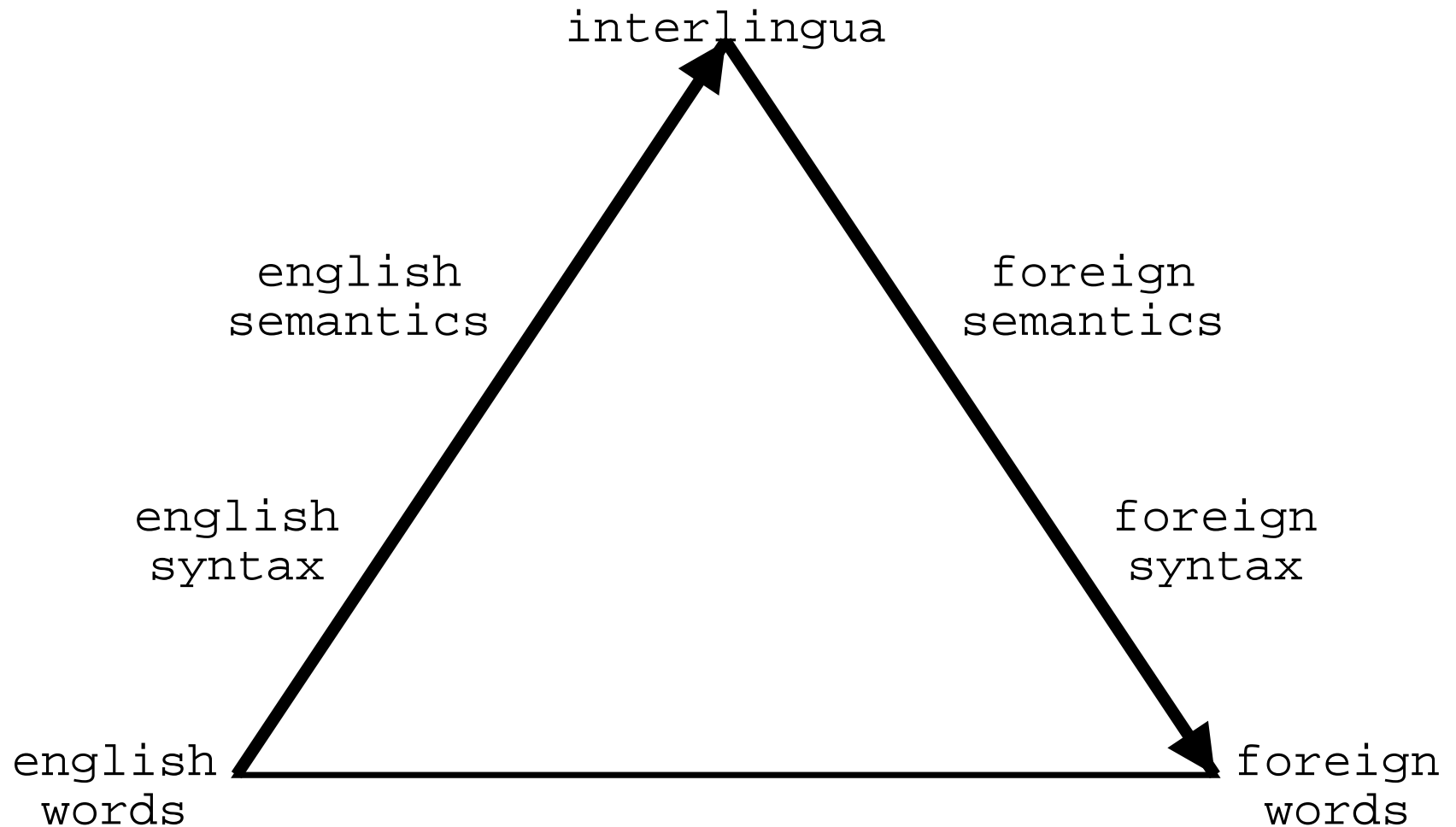
Translation Model

- Goal of the Translation Model:
Match foreign input to English output

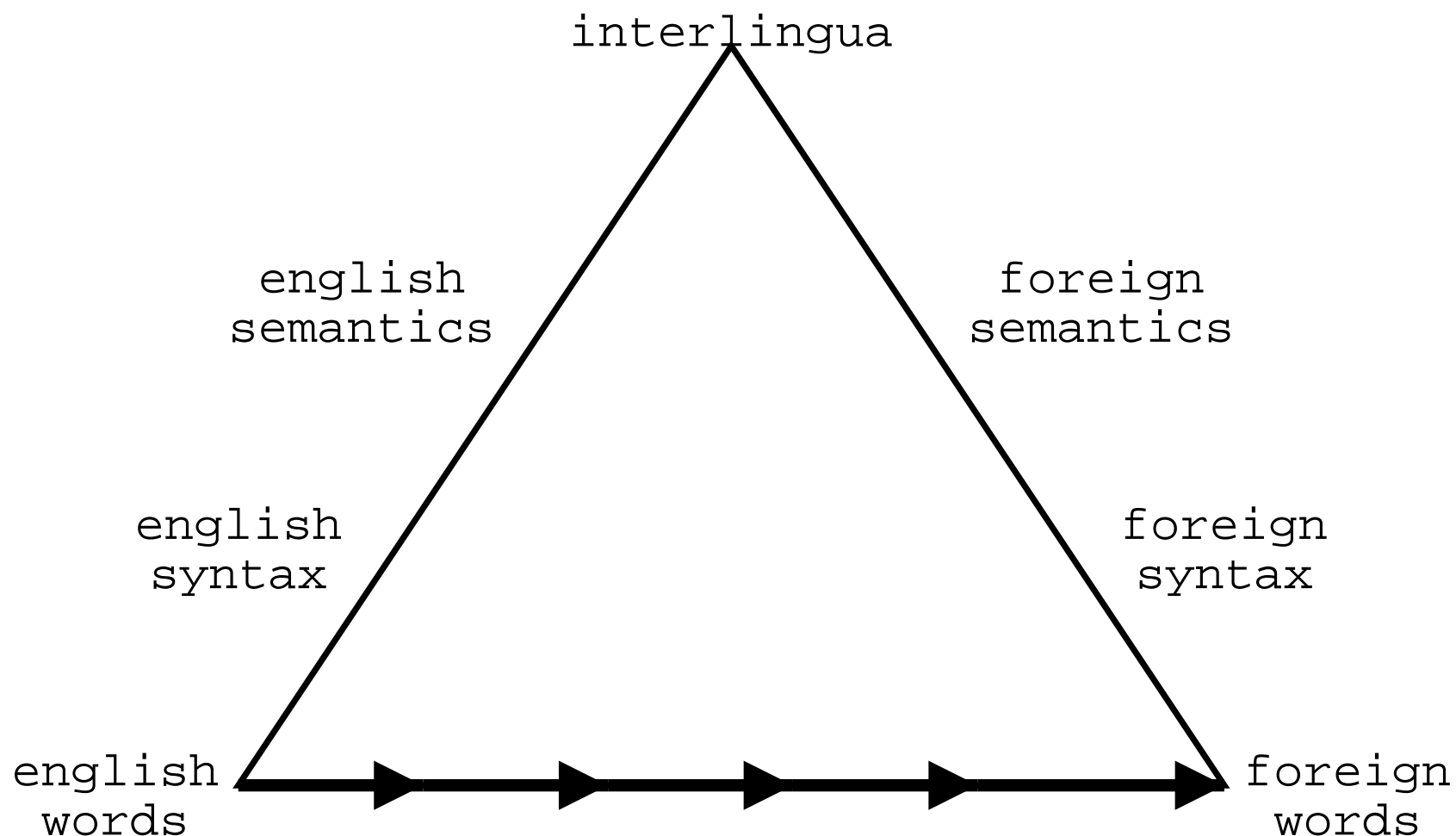
Overview: Translation Model

- **Machine translation pyramid**
- Statistical modeling and IBM Model 4
- EM algorithm
- Word alignment
- Flaws of word-based translation
- Phrase-based translation
- Syntax-based translation

The Machine Translation Pyramid



The Machine Translation Pyramid



however, the currently best performing statistical machine translation systems are still crawling at the bottom.

Overview: Translation Model

- Machine translation pyramid
- **Statistical modeling and IBM Model 4**
- EM algorithm
- Word alignment
- Flaws of word-based translation
- Phrase-based translation
- Syntax-based translation

Statistical Modeling

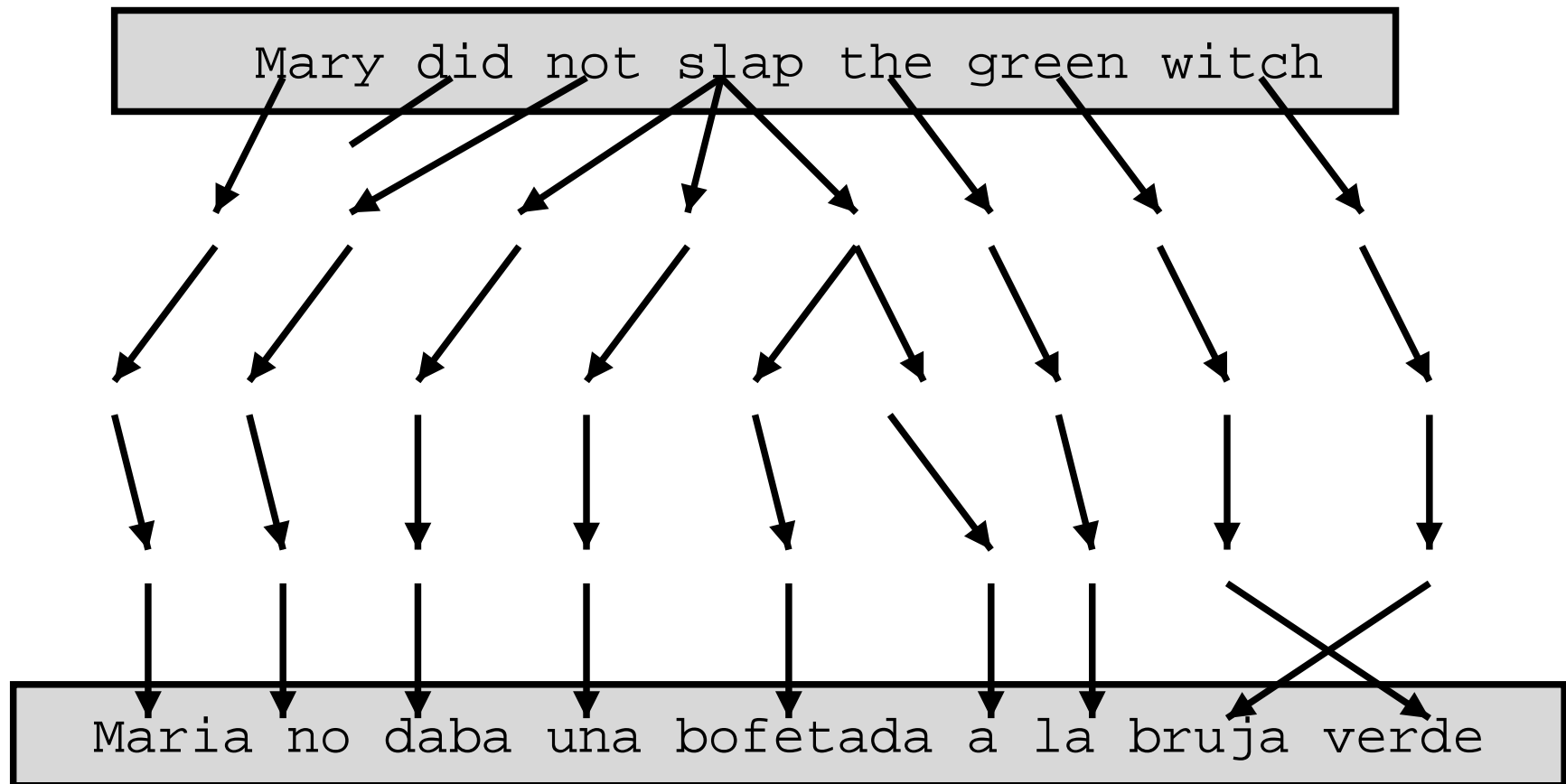
Mary did not slap the green witch



Maria no daba una bofetada a la bruja verde

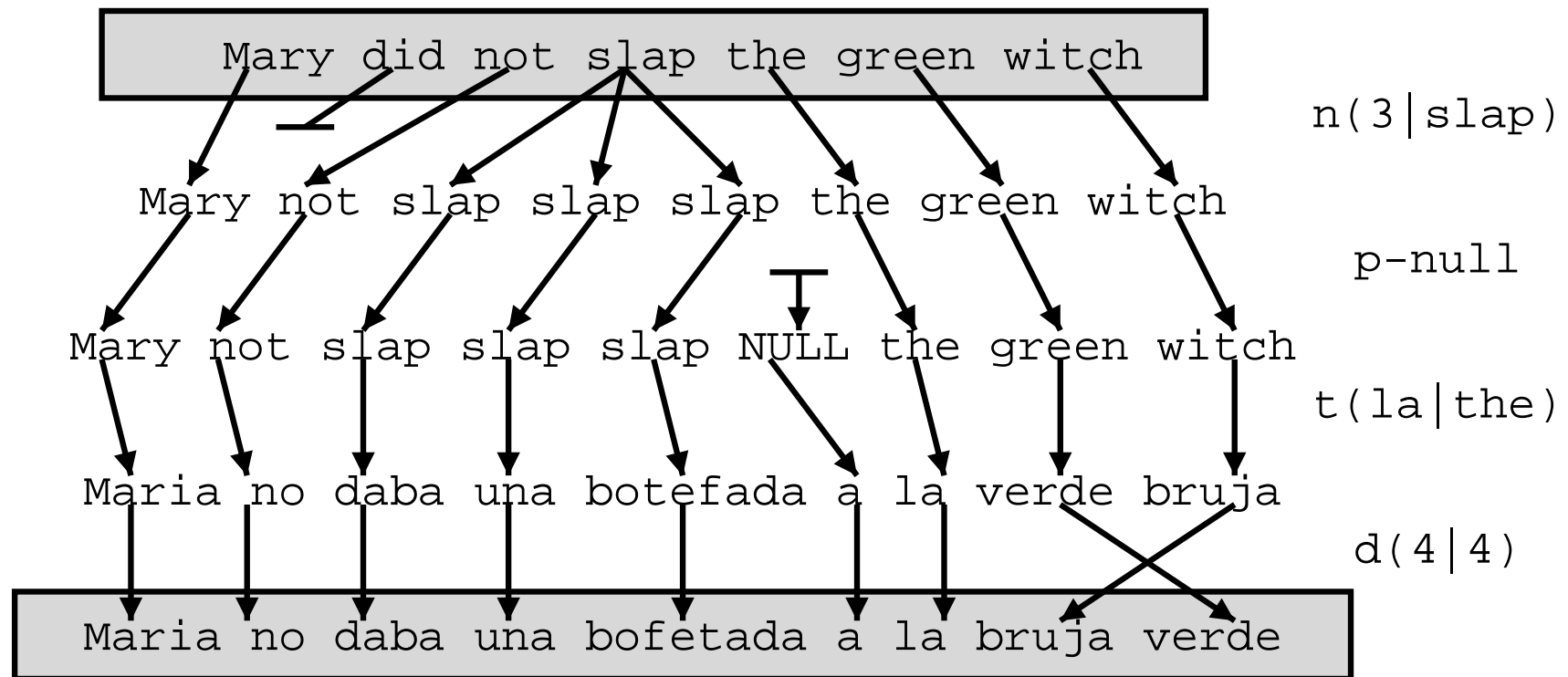
- Learn $P(f|e)$ from a Parallel Corpus
- Not Sufficient Data to Estimate $P(f|e)$ Directly

Statistical Modeling (2)



- Break the Process into Smaller Steps

Statistical Modeling (3)



- Probabilities for Smaller Steps can be Learned

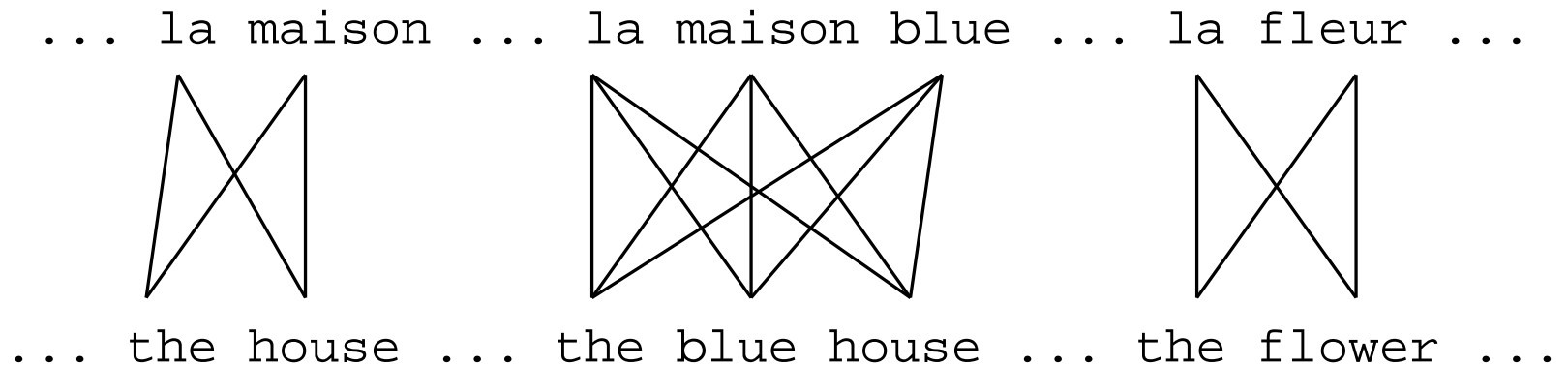
Statistical Modeling (4)

- Generate a Story How an English String e Gets to be a Foreign String f
 - Choices in Story are Decided by Reference to Parameters
 - e.g., $p(\text{bruja}|\text{witch})$
- Formula for $P(f|e)$ in Terms of Parameters
 - usually long and hairy, but mechanical to extract from the story
- Training to Obtain Parameter Estimates from Possibly Incomplete Data
 - off-the-shelf EM

Overview: Translation Model

- Machine translation pyramid
- Statistical modeling and IBM Model 4
- **EM algorithm**
- Word alignment
- Flaws of word-based translation
- Phrase-based translation
- Syntax-based translation

Parallel Corpora



- Incomplete Data

- English and foreign words, but no connections between them

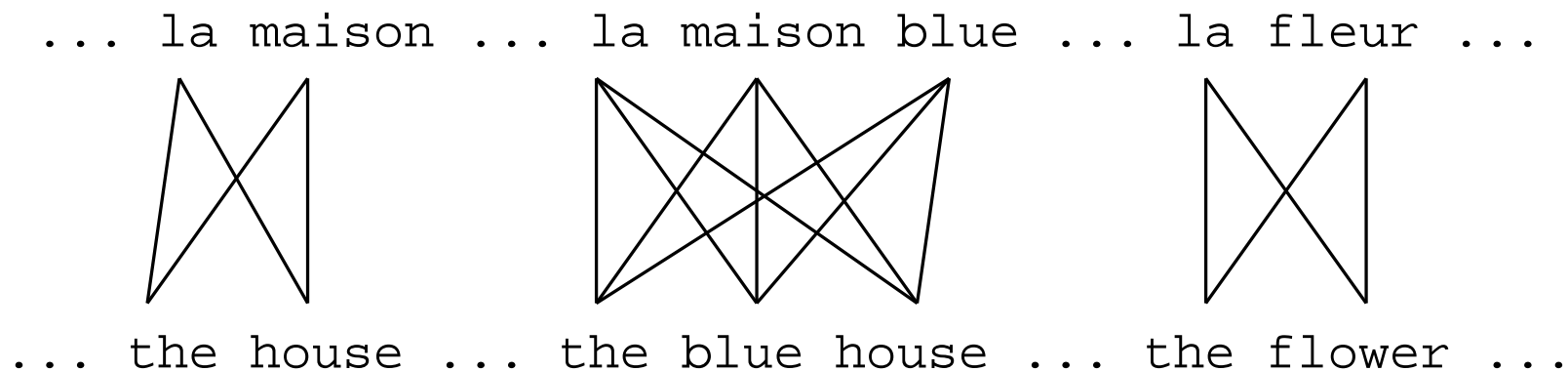
- Chicken and Egg Problem

- if we had the connections, we could estimate the parameters of our generative story
- if we had the parameters, we could estimate the connections

EM Algorithm

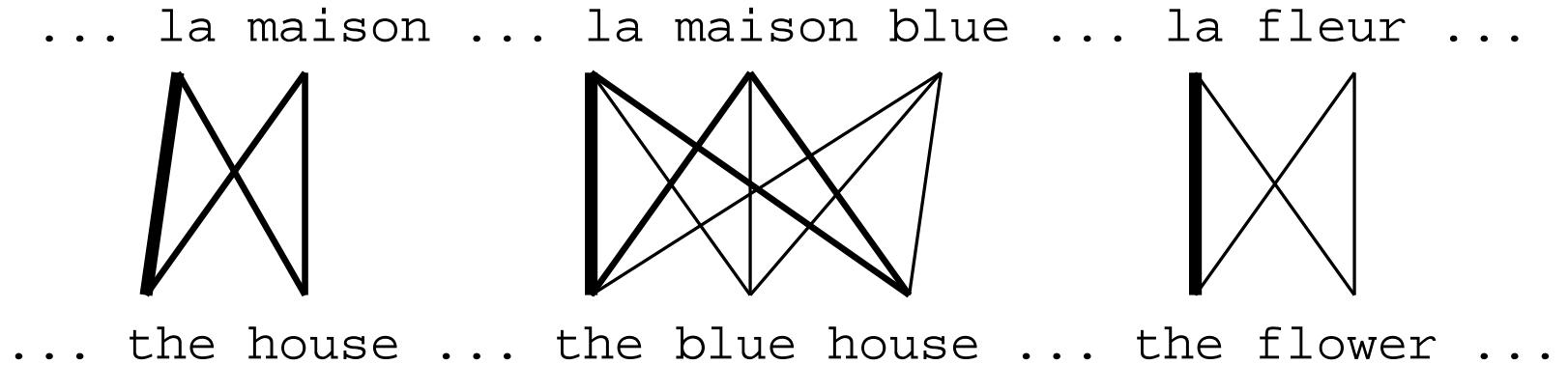
- Incomplete Data
 - if we had complete data, would could estimate model
 - if we had model, we could fill in the gaps in the data
- EM in a Nutshell
 - initialize model parameters (e.g. uniform)
 - assign probabilities to the missing data
 - estimate model parameters from completed data
 - iterate

EM Algorithm (2)



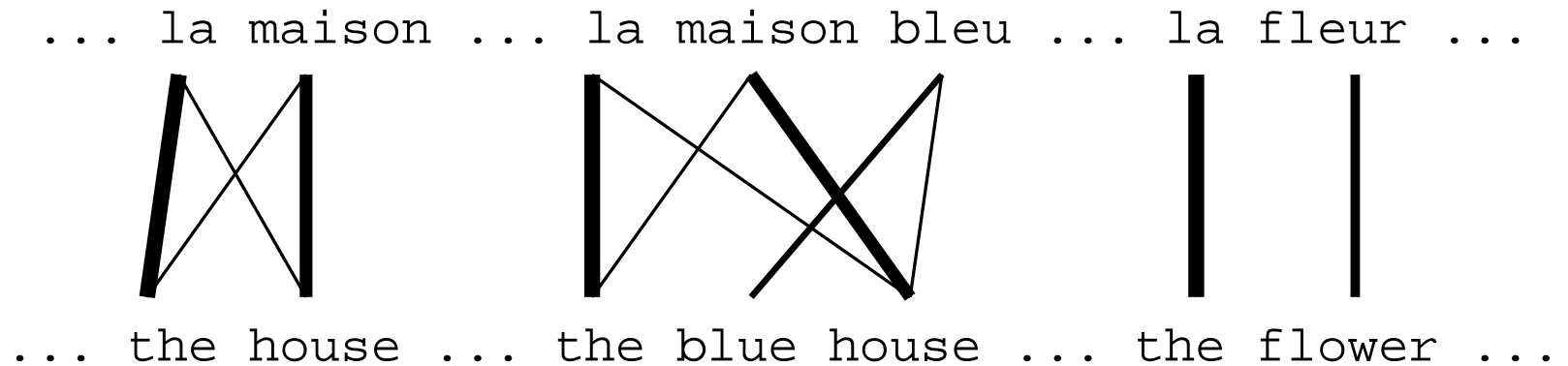
- Initial Step: all Connections Equally Likely
- Model Learns that, e.g., **la** is Often Connected with **the**

EM Algorithm (3)



- After One Iteration
- Connections, e.g., between **la** and **the** are More Likely

EM Algorithm (4)



- After Another Iteration
- It Becomes Apparent that Connections, e.g., between **fleur** and **flower** are More Likely (Pigeon Hole Principle)

EM Algorithm (5)

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

- Convergence
- Inherent Hidden Structure Revealed by EM

EM Algorithm (6)

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



$p(\text{la}|\text{the}) = 0.453$
 $p(\text{le}|\text{the}) = 0.334$
 $p(\text{maison}|\text{house}) = 0.876$
 $p(\text{bleu}|\text{blue}) = 0.563$
...

- Parameter Estimation from the Connected Corpus

More detail on the IBM Models

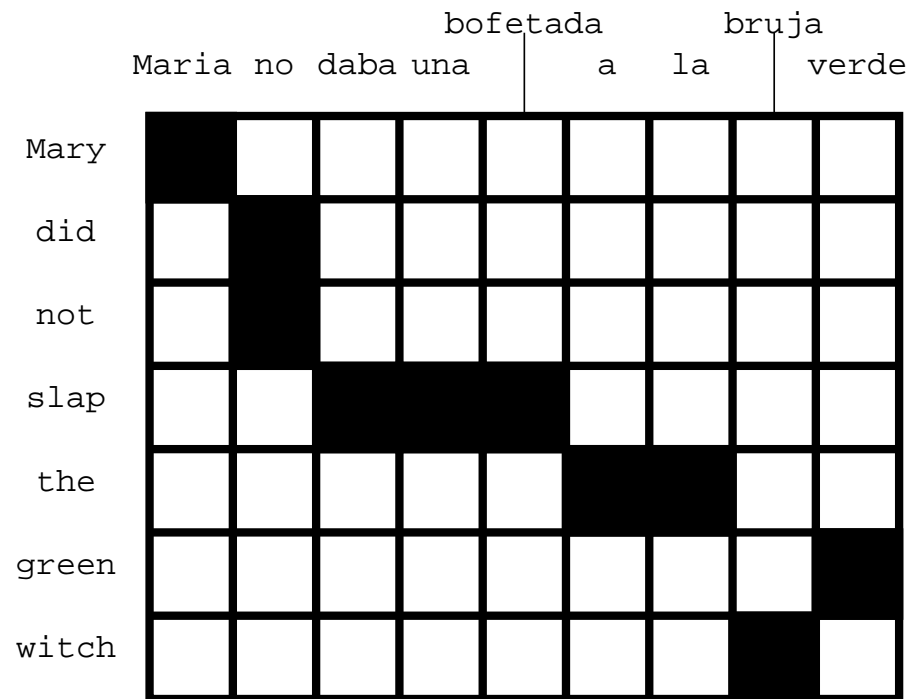
- “A Statistical MT Tutorial Workbook” (Knight, 1999)
- “The Mathematics of Statistical Machine Translation” (Brown et al., 1993)
- Downloadable Software: Giza++, ReWrite Decoder

Overview: Translation Model

- Machine translation pyramid
- Statistical modeling and IBM Model 4
- EM algorithm
- **Word alignment**
- Flaws of word-based translation
- Phrase-based translation
- Syntax-based translation

Word Alignment

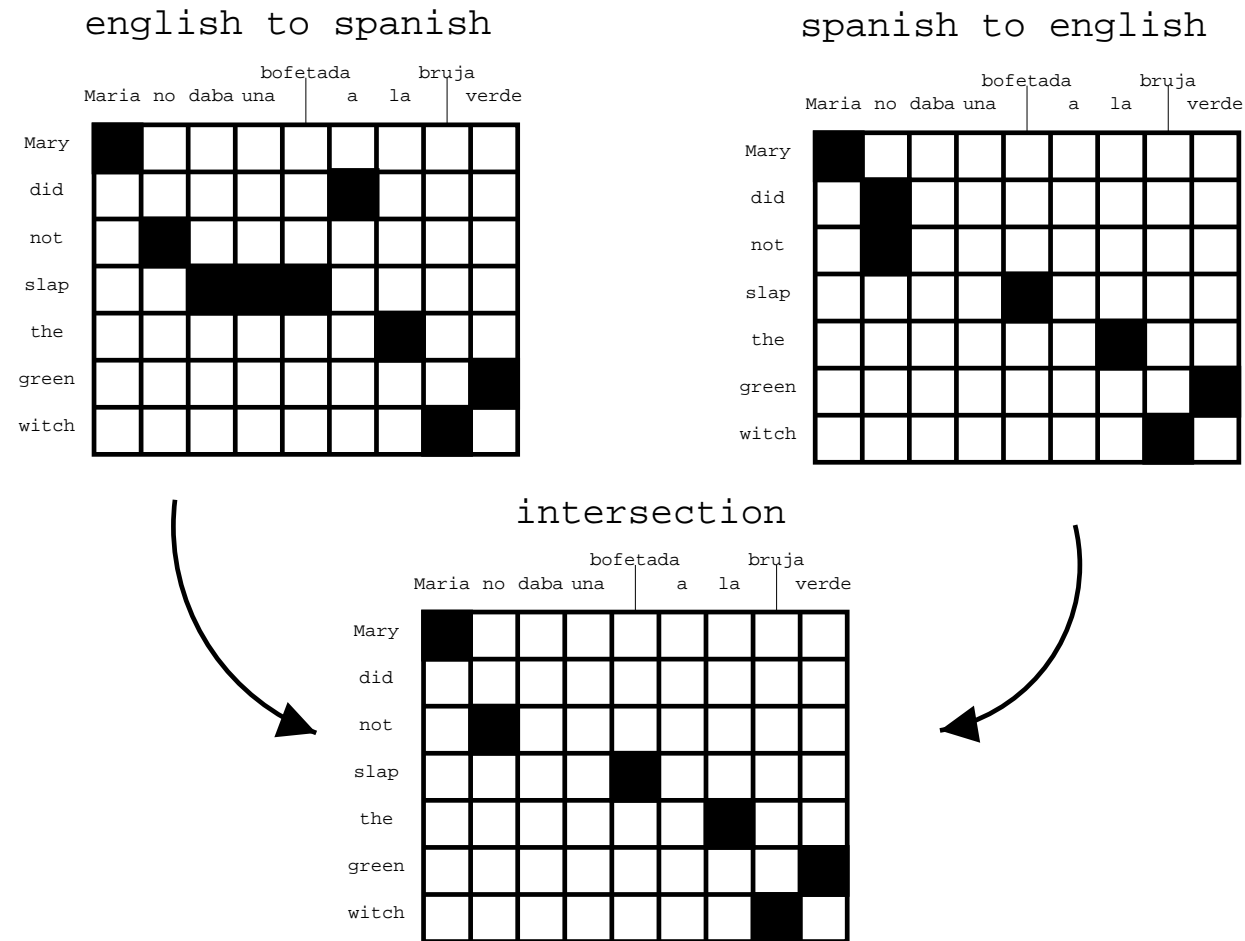
- Notion of Word Alignments Valuable
- Trained Humans can Achieve High Agreement
- Shared Task at Data-Driven MT Workshop at NAACL/HLT



Improved Word Alignments

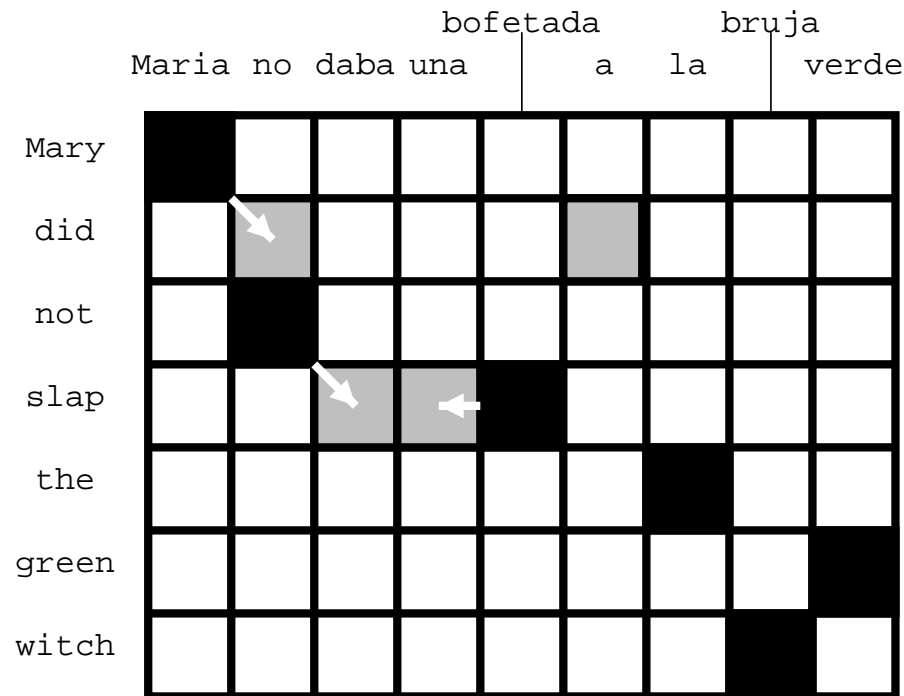
- Improving IBM Model Word Alignments with Heuristics [Och and Ney, 2000, Koehn et al., 2003]
 - one-to-many problem of IBM Models
 - bidirectionally aligned corpora $e \rightarrow f, f \rightarrow e$
 - take intersection of alignment points
(high precision, low recall)
 - grow additional alignment points
(increase recall while preserving precision)

Improved Word Alignments (2)



- Intersection of Bidirectional Alignments

Improved Word Alignments (3)



- Grow Additional Alignment Points

Improved Word Alignments (4)

- Heuristics for Adding Alignment Points

- only to directly neighboring
- also to diagonally neighboring
- also to non-neighboring
- prefer English-foreign or foreign-to-English
- use lexical probabilities or frequencies
- extend only to unaligned words
- ...

⇒ No Clear Advantage to any Strategy

- depends on corpus size
- depends on language pair

Overview: Translation Model

- Machine translation pyramid
- Statistical modeling and IBM Model 4
- EM algorithm
- Word alignment
- **Flaws of word-based translation**
- Phrase-based translation
- Syntax-based translation

Flaws of Word-Based MT

- Multiple English Words for one German Word

German:	Zeitmangel	erschwert	das	Problem	.
Gloss:	LACK OF TIME	MAKES MORE DIFFICULT	THE	PROBLEM	.
Correct translation:	Lack of time makes the problem more difficult.				
MT output:	Time makes the problem .				

- Phrasal Translation

German:	Eine	Diskussion	erübrigt	sich	demnach
Gloss:	A	DISCUSSION	IS MADE UNNECESSARY	ITSELF	THEREFORE
Correct translation:	Therefore, there is no point in a discussion.				
MT output:	A debate turned therefore .				

Flaws of Word-Based MT (2)

- Syntactic Transformations

German: Das ist der Sache nicht angemessen .

Gloss: THAT IS THE MATTER NOT APPROPRIATE .

Correct translation: That is not appropriate for this matter .

MT output: That is the thing is not appropriate .

German: Den Vorschlag lehnt die Kommission ab .

Gloss: THE PROPOSAL REJECTS THE COMMISSION OFF .

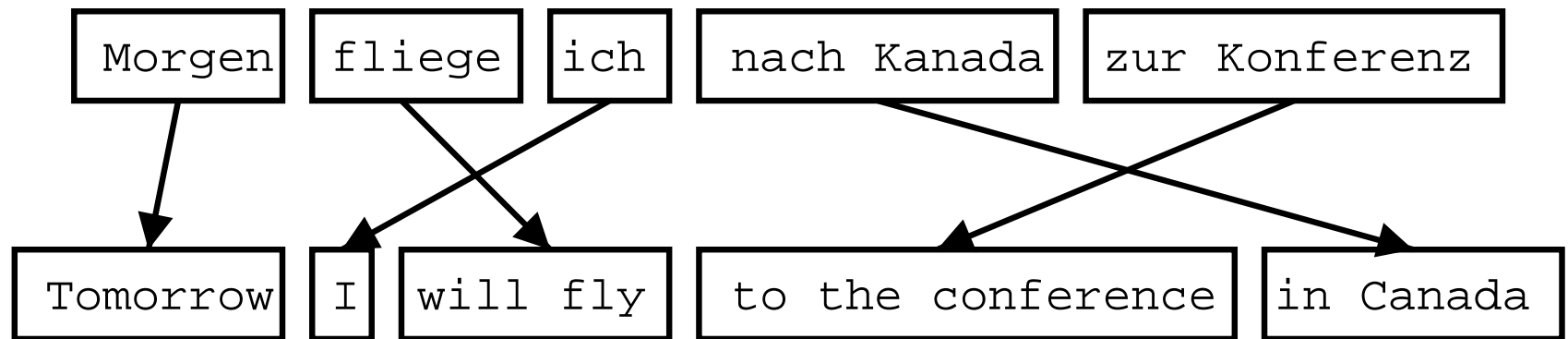
Correct translation: The commission rejects the proposal .

MT output: The proposal rejects the commission .

Overview: Translation Model

- Machine translation pyramid
- Statistical modeling and IBM Model 4
- EM algorithm
- Word alignment
- Flaws of word-based translation
- **Phrase-based translation**
- Syntax-based translation

Phrase-Based Translation



- Foreign Input is Segmented in Phrases
 - any sequence of words, not necessarily linguistically motivated
- Each Phrase is Translated into English
- Phrases are Reordered

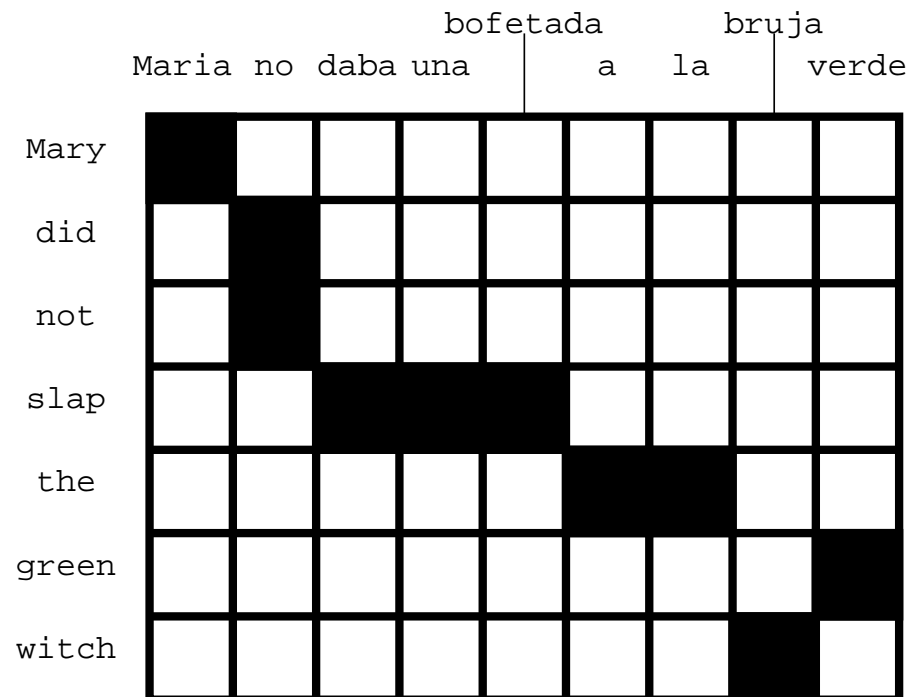
Advantages of Phrase-Based Translation

- Many-to-Many Translation
- Use of Local Context in Translation
- Allows Translation of Non-Compositional Phrases
- The More Data, the Longer Phrases can be Learned

Three Phrase-Based Translation Models

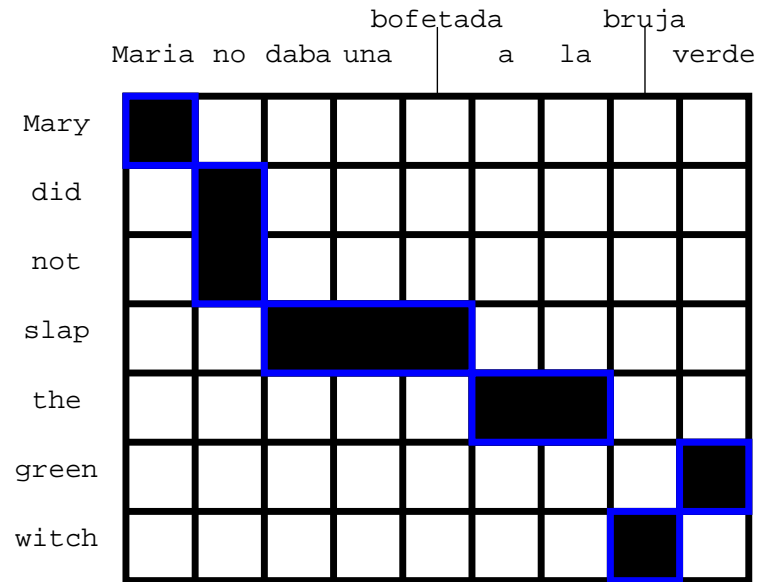
- Word Alignment Induced Phrase Model
[Koehn et al., 2003]
- Alignment Templates [Och et al., 1999]
- Joint Phrase Model [Marcu and Wong, 2002]

Word Alignment Induced Phrases



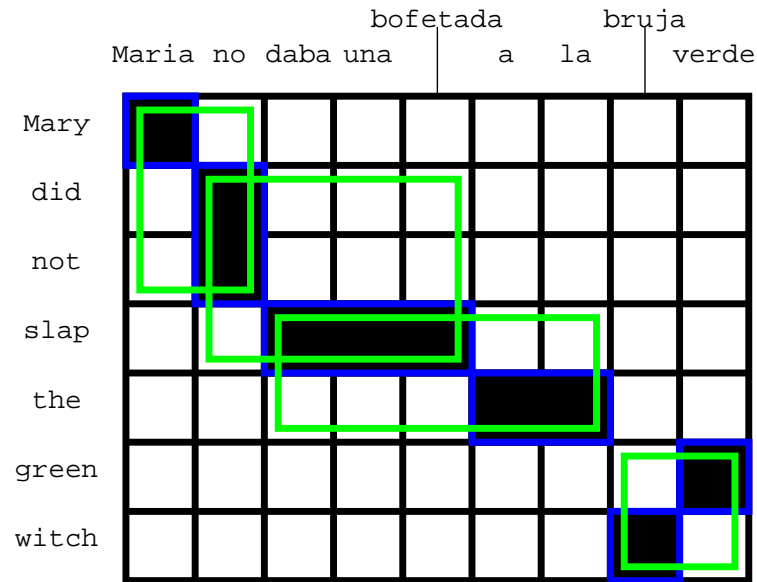
- Collect All Phrase Pairs that are Consistent with the Word Alignment
 - a phrase alignment has to contain all alignment points for all words it covers

Word Alignment Induced Phrases (2)



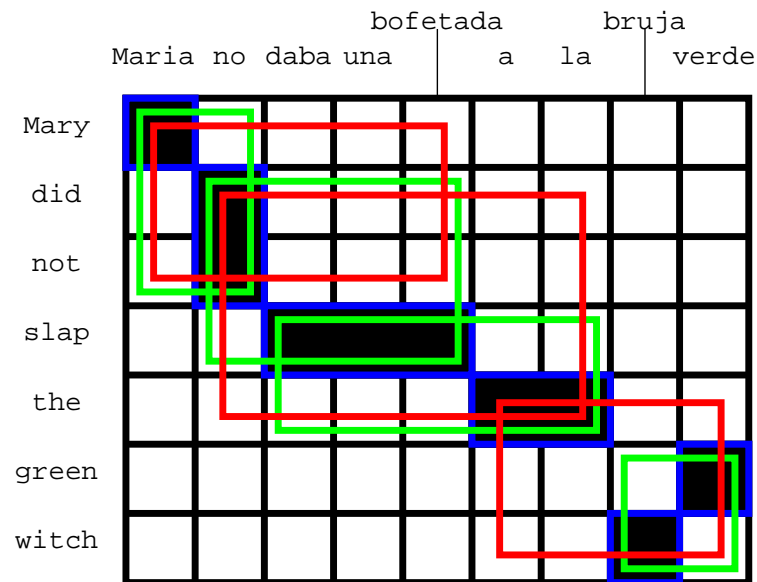
(Maria, Mary), (no, did not), (slap, daba una bofetada), (a la, the), (bruja, witch), (verde, green)

Word Alignment Induced Phrases (3)



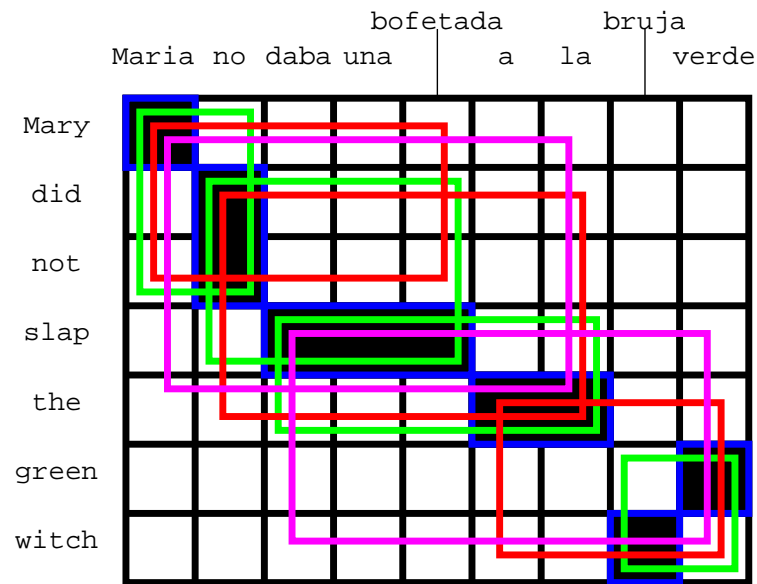
(Maria, Mary), (no, did not), (slap, daba una bofetada), (a la, the), (bruja, witch),
 (verde, green), (Maria no, Mary did not), (no daba una bofetada, did not slap),
 (daba una bofetada a la, slap the), (bruja verde, green witch)

Word Alignment Induced Phrases (4)



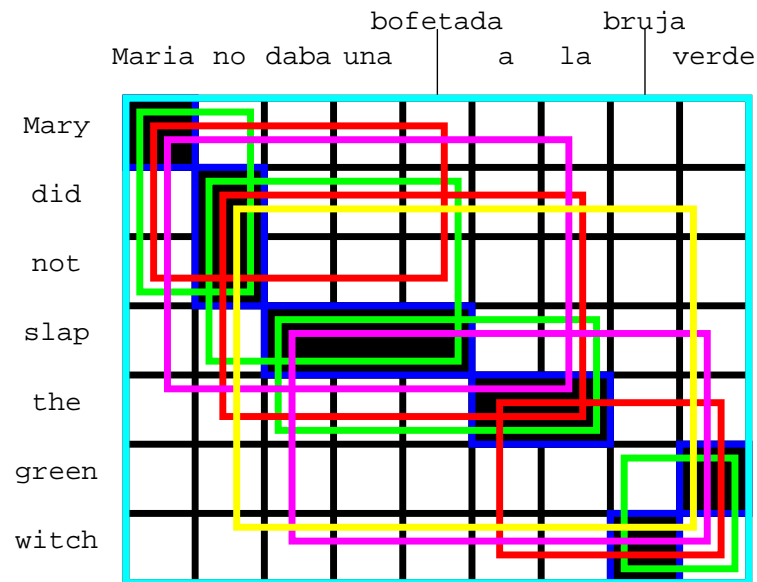
(Maria, Mary), (no, did not), (slap, daba una bofetada), (a la, the), (bruja, witch),
 (verde, green), (Maria no, Mary did not), (no daba una bofetada, did not slap),
 (daba una bofetada a la, slap the), (bruja verde, green witch),
 (Maria no daba una bofetada, Mary did not slap),
 (no daba una bofetada a la, did not slap the), (a la bruja verde, the green witch)

Word Alignment Induced Phrases (5)



(Maria, Mary), (no, did not), (slap, daba una bofetada), (a la, the), (bruja, witch),
 (verde, green), (Maria no, Mary did not), (no daba una bofetada, did not slap),
 (daba una bofetada a la, slap the), (bruja verde, green witch),
 (Maria no daba una bofetada, Mary did not slap),
 (no daba una bofetada a la, did not slap the), (a la bruja verde, the green witch),
 (Maria no daba una bofetada a la, Mary did not slap the),
 (daba una bofetada a la bruja verde, slap the green witch)

Word Alignment Induced Phrases (6)



- (Maria, Mary), (no, did not), (slap, daba una bofetada), (a la, the), (bruja, witch),
- (verde, green), (Maria no, Mary did not), (no daba una bofetada, did not slap),
- (daba una bofetada a la, slap the), (bruja verde, green witch),
- (Maria no daba una bofetada, Mary did not slap),
- (no daba una bofetada a la, did not slap the), (a la bruja verde, the green witch),
- (Maria no daba una bofetada a la, Mary did not slap the),
- (daba una bofetada a la bruja verde, slap the green witch),
- (no daba una bofetada a la bruja verde, did not slap the green witch),
- (Maria no daba una bofetada a la bruja verde, Mary did not slap the green witch)

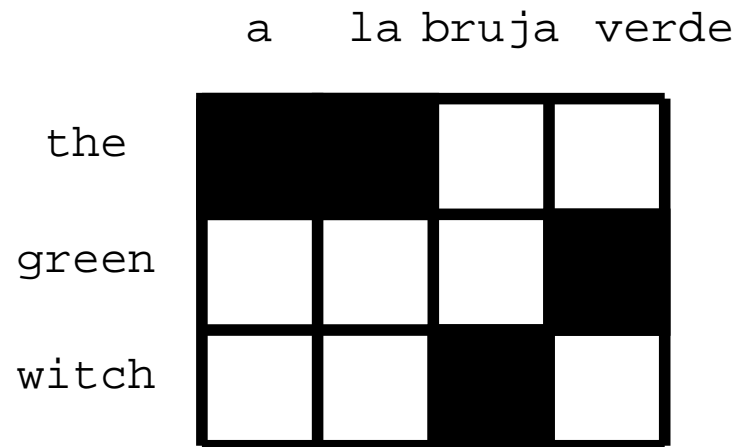
Word Alignment Induced Phrases (7)

- Given the Collected Phrase Pairs,
Estimate the Phrase Translation Probability Distribution
by Relative Frequency:

$$\phi(\bar{f}|\bar{e}) = \frac{\text{count}(\bar{f},\bar{e})}{\sum_{\bar{f}} \text{count}(\bar{f},\bar{e})}$$

- No Smoothing is Performed

Word Alignment Induced Phrases (8)



- Lexical Weighting:

$$p_w(\bar{f}|\bar{e}, a) = \prod_{i=1}^n \frac{1}{|\{j|(i,j) \in a\}|} \sum_{\forall (i,j) \in a} w(f_i|e_j)$$

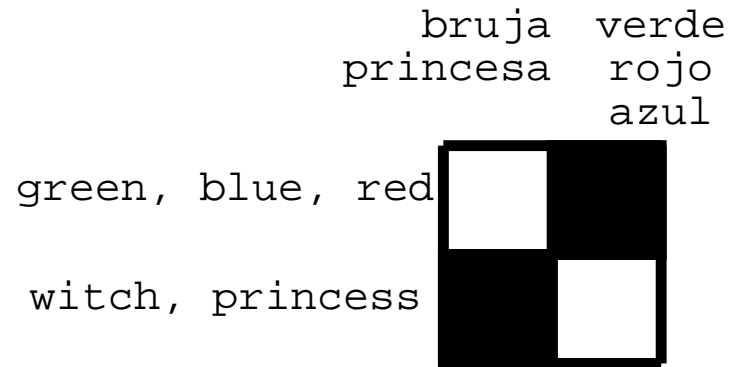
$$p_w(\text{a la bruja verde}|\text{the green witch}) =$$

$$w(\text{a}|\text{the}) \times w(\text{la}|\text{the}) \times$$

$$w(\text{verde}|\text{green}) \times$$

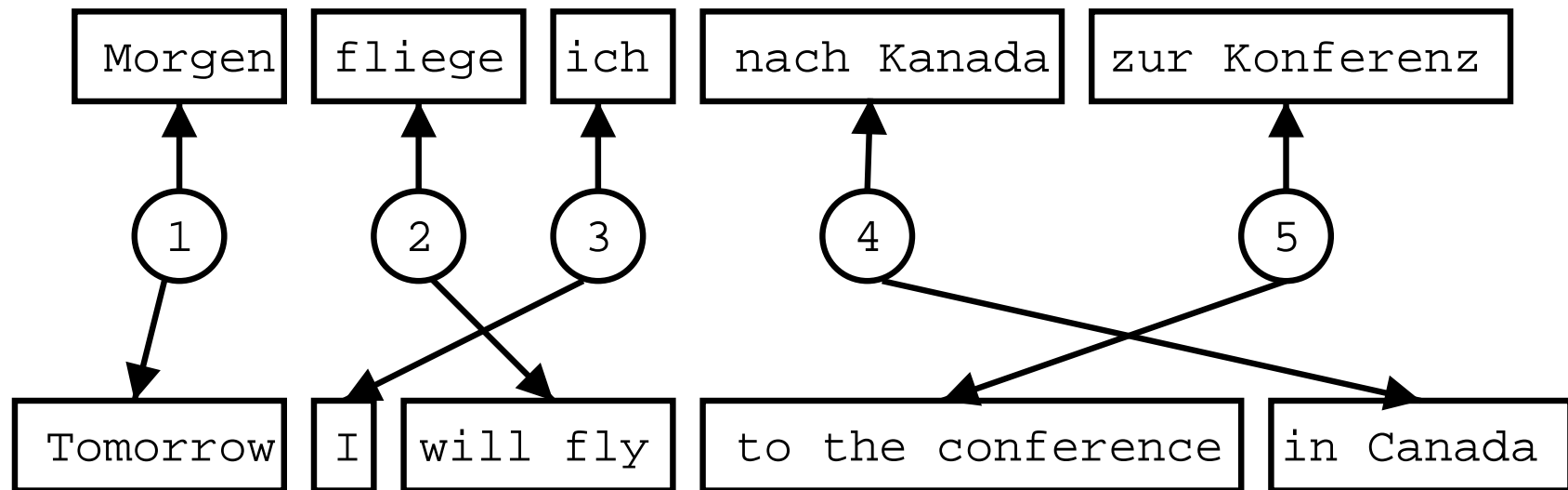
$$w(\text{bruja}|\text{witch})$$

Alignment Templates [Och et al., 1999]



- Word Classes instead of Words
 - alignment templates instead of phrases
 - more reliable statistics for translation table
 - smaller translation table
 - more complex decoding
- Same Lexical Weighting

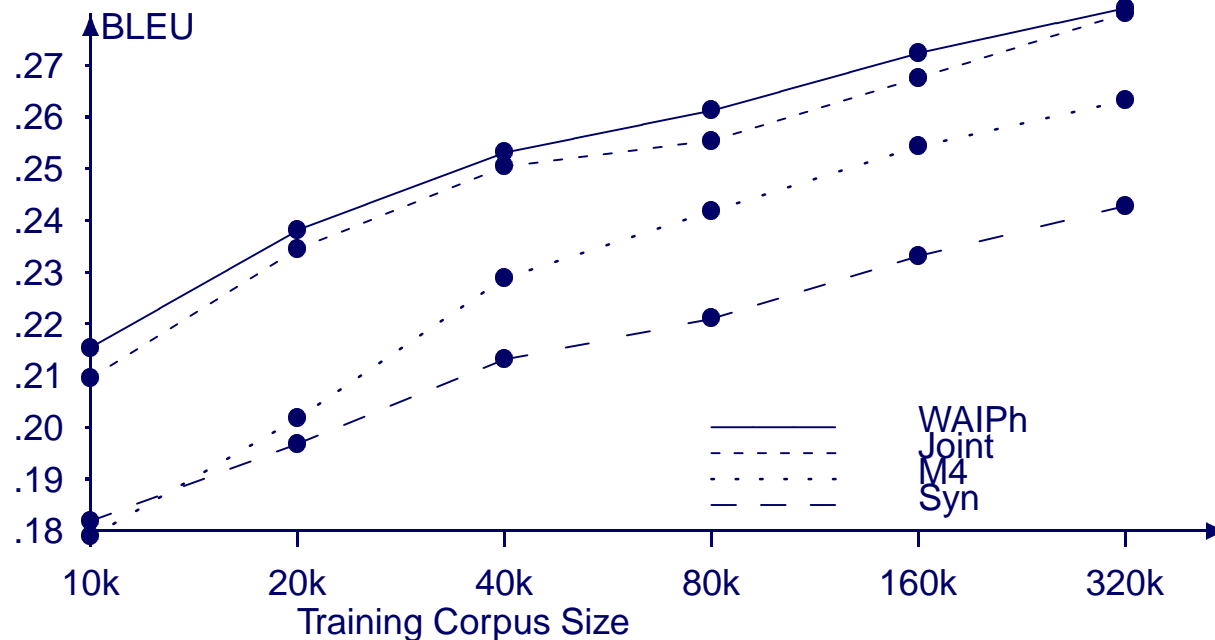
Joint Phrase Model



- Direct Phrase Alignment of Parallel Corpus
[Marcu and Wong, 2002]
- Generative Story
 - a number of concepts are created
 - each concept generates a foreign and English phrase
 - the English phrases are reordered

Evaluation of Phrase Models

- Direct Comparison of Models [Koehn et al., 2003]
 - results improve log-linear with training corpus size
 - WAIPh slightly better than Joint (same decoder, same LM)
 - better than IBM Model 4 (different decoder)
 - using only phrases that are syntactic constituents hurts



Evaluation of Phrase Models (2)

- Different Language Pairs
 - results for WAIPh
 - better than IBM Model 4
 - lexical weighting always helps

Language Pair	Model4	Phrase	Lex
English-German	0.20	0.24	0.24
French-English	0.28	0.33	0.34
English-French	0.26	0.31	0.32
Finnish-English	0.22	0.27	0.28
Swedish-English	0.31	0.35	0.36
Chinese-English	0.12	0.14	0.14

Limits of Phrase Models

- Non-Contiguous Phrases

- German: *Ich habe das Auto gekauft*
- English: *I bought the car*
- good phrase pair: *habe ... gekauft* == *bought*

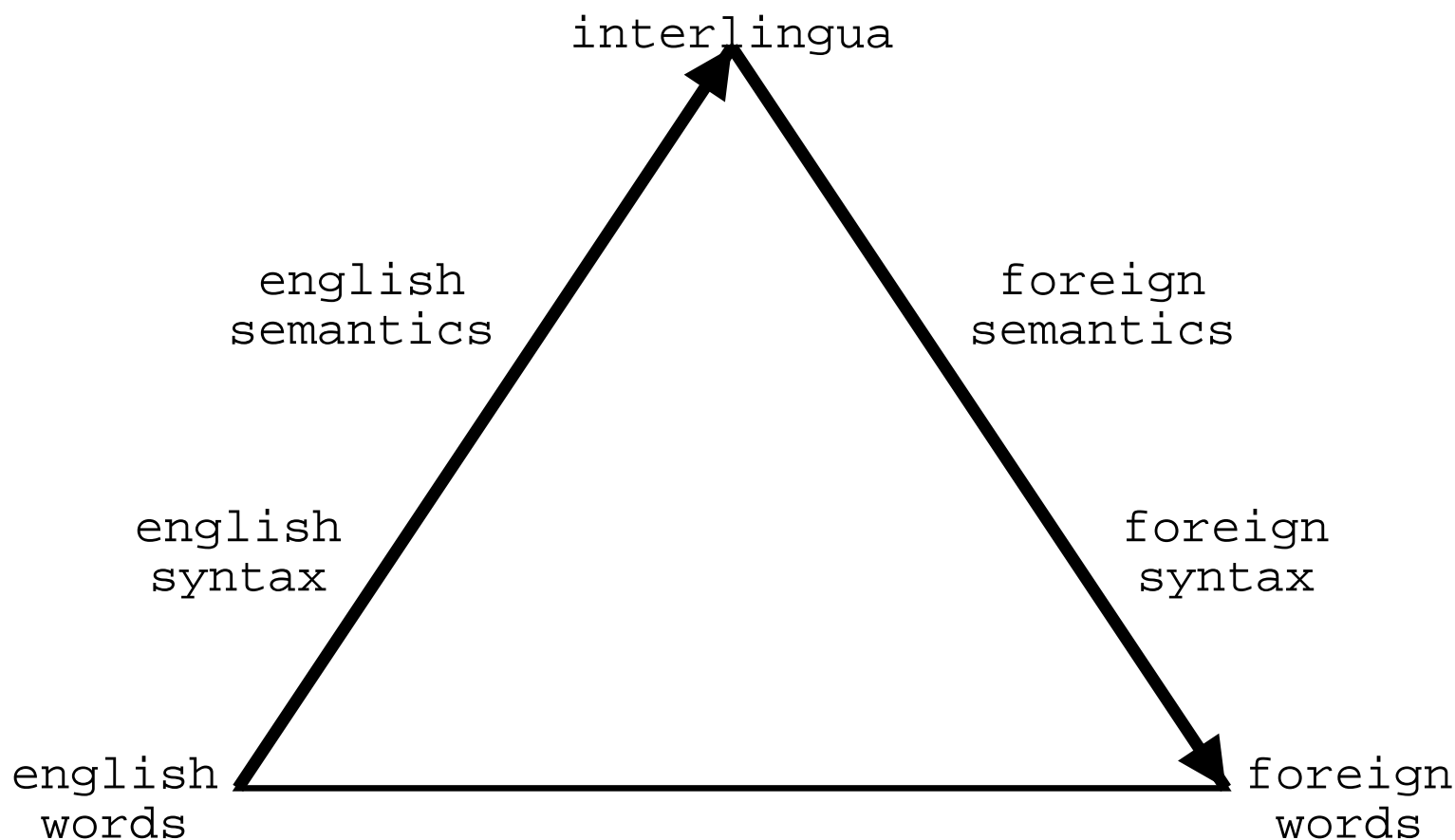
- Syntactic Transformations

- German: *Den Antrag verabschiedet das Parlament*
- English gloss: *The draft approves the Parliament*
- case marking that indicates that “the draft” is object is lost during translation

Overview: Translation Model

- Machine translation pyramid
- Statistical modeling and IBM Model 4
- EM algorithm
- Word alignment
- Flaws of word-based translation
- Phrase-based translation
- **Syntax-based translation**

Syntax-Based Translation

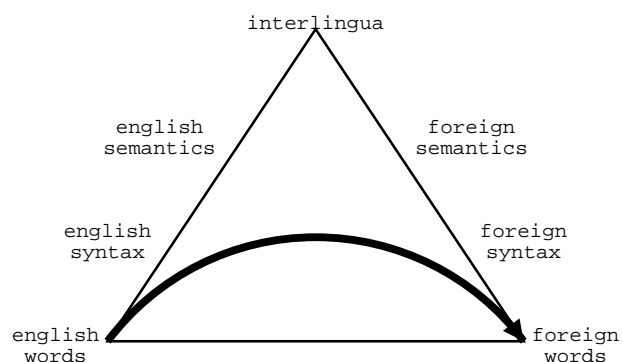


- Remember the Pyramid

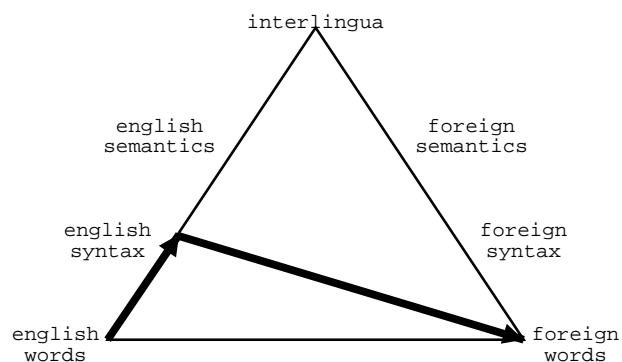
Advantages of Syntax-Based Translation

- Reordering for Syntactic Reasons
 - e.g., move German object to end of sentence
- Better Explanation for Function Words
 - e.g., prepositions, determiners
- Conditioning to Syntactically Related Words
 - translation of verb may depend on subject or object
- Use of Syntactic Language Models

Syntax-Based Translation Models



Wu [1997], Alshawi et al. [1998]



Yamada and Knight [2001]

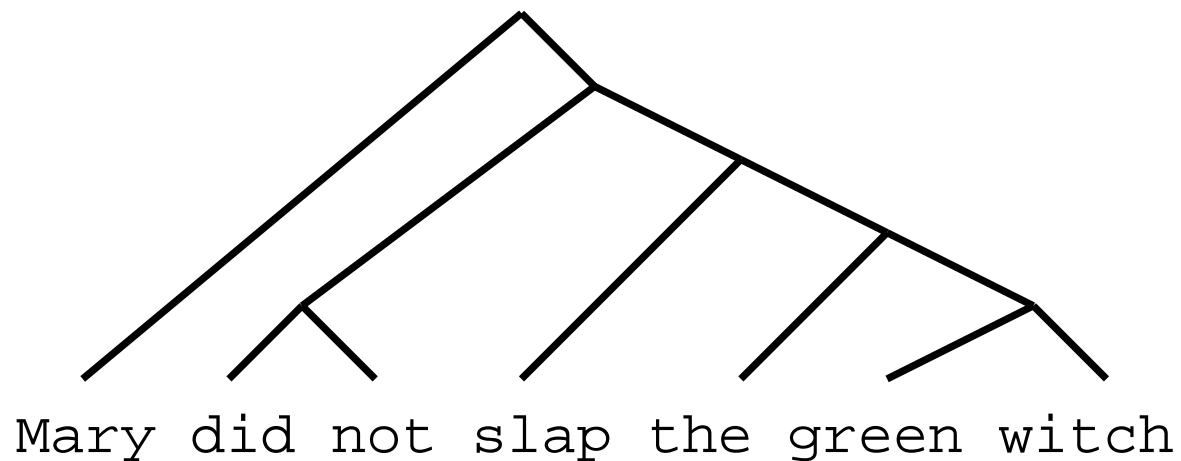
Inversion Transduction Grammars

- Generation of both English and Foreign Trees [Wu, 1997]
- Rules (Binary and Unary)
 - $A \rightarrow A_1 A_2 | A_1 A_2$
 - $A \rightarrow A_1 A_2 | A_2 A_1$
 - $A \rightarrow e | f$
 - $A \rightarrow e | *$
 - $A \rightarrow * | f$

⇒ Common Binary Tree Required

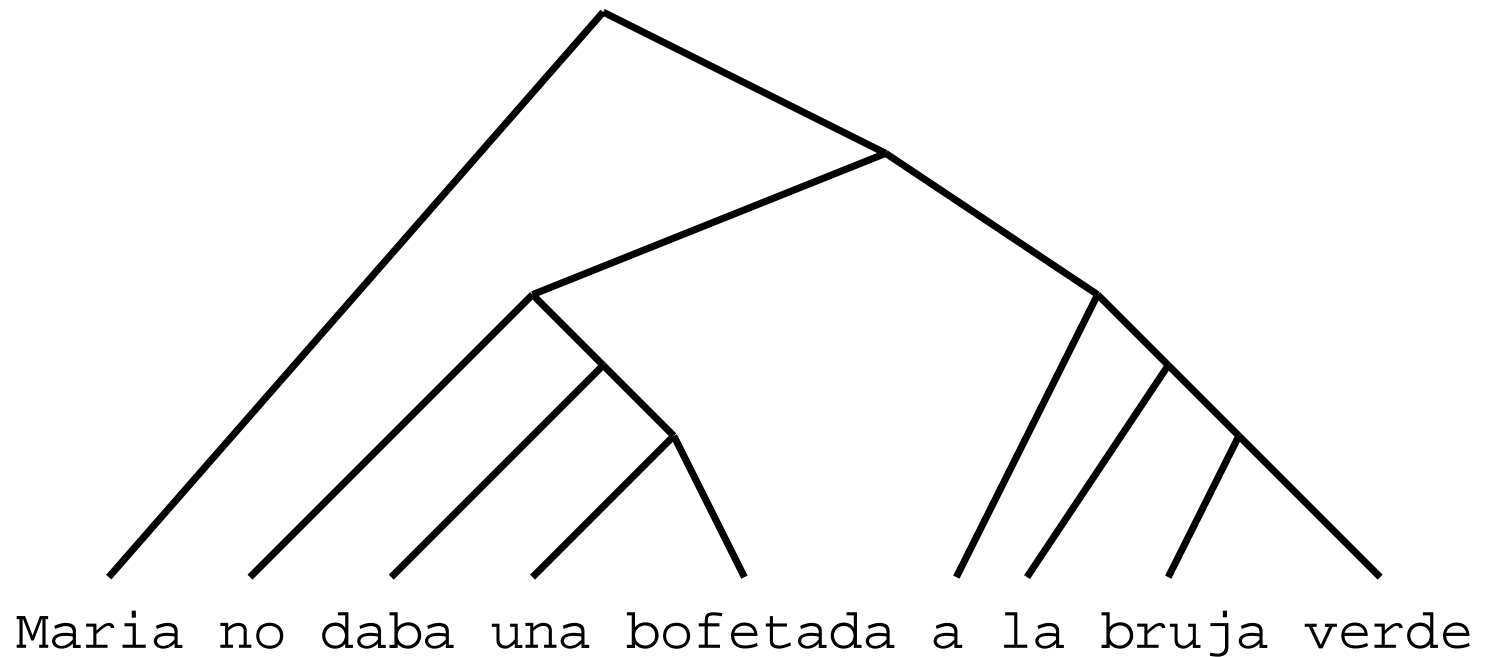
- limits the complexity of reorderings

Syntax Trees



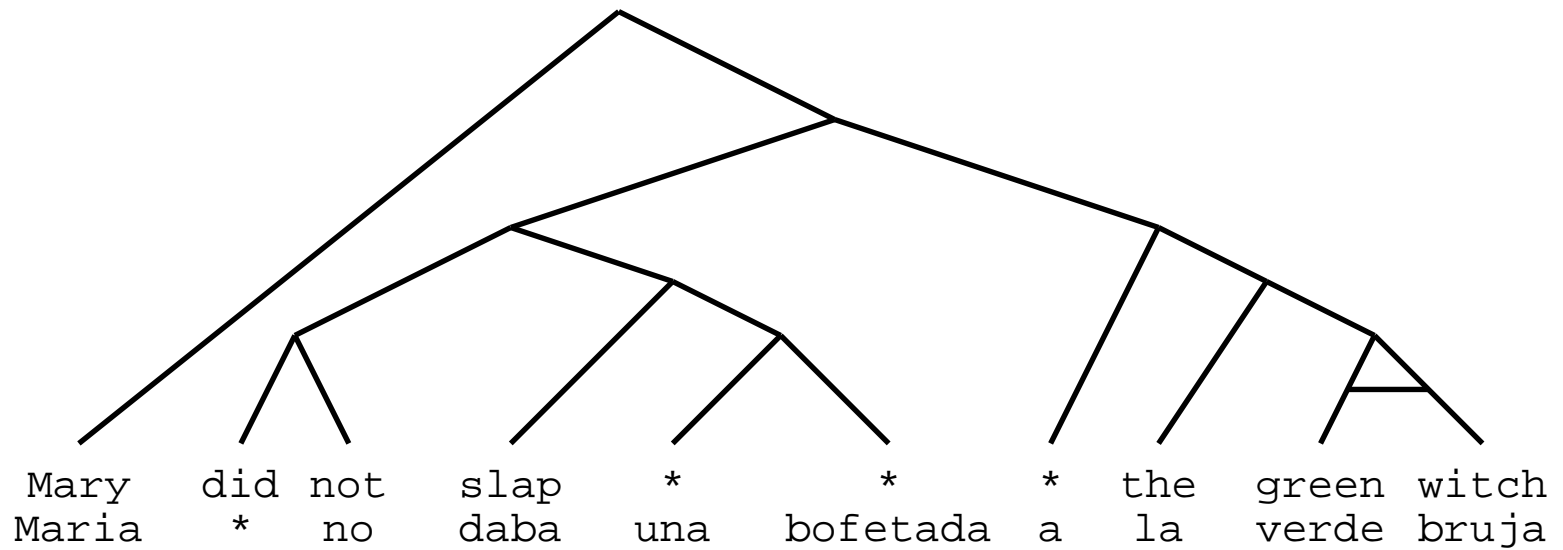
- English Binary Tree

Syntax Trees (2)



- Spanish Binary Tree

Syntax Trees (3)

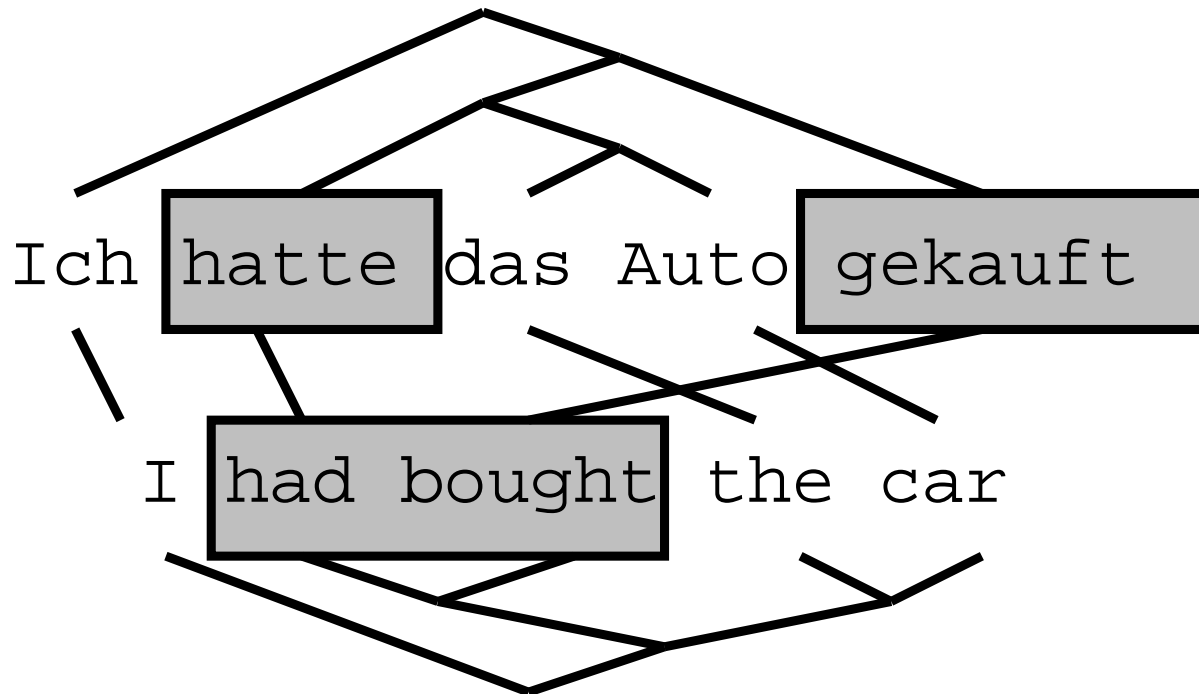


- Combined Tree with Reordering of Spanish

Hierarchical Transduction Models

- Based on Finite State Transducers [Alshawi et al., 1998]
 - also common binary tree required
 - lexicalized non-terminal rules
- Generation of Sentence Pair
 1. create initial **head word** (e.g., [daba : slap])
 2. extend head word by **adding dependents** (e.g., [bruja : witch]);
foreign and English could be placed on different sides of head;
dependents could be single word, empty, or phrases
 3. pick one of the dependents as new head word for extension (step 2);
or terminate

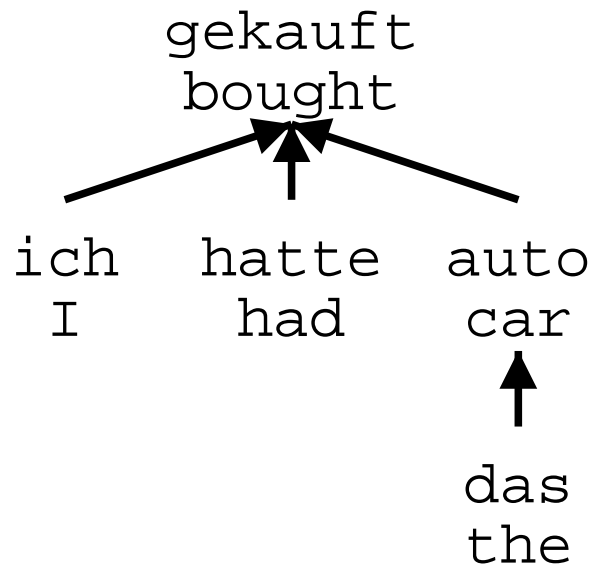
Common Binary Tree Requirement



- No Common Binary Tree Possible
- Maybe Languages are Syntactically too Different?

⇒ Jump Ahead to Semantics

Dependency Structure

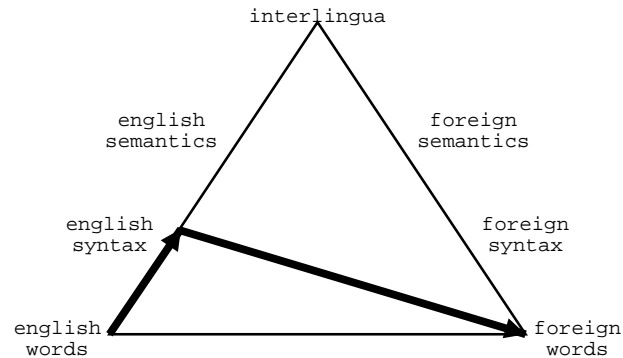


- Common Dependency Tree
- Interest in Dependency-Based Translation Models
 - e.g. Czech-English [Cmejrek et al., 2003]
 - current systems mixed statistical/rule-based
 - probably good generation system necessary

Direct Correspondence Assumption

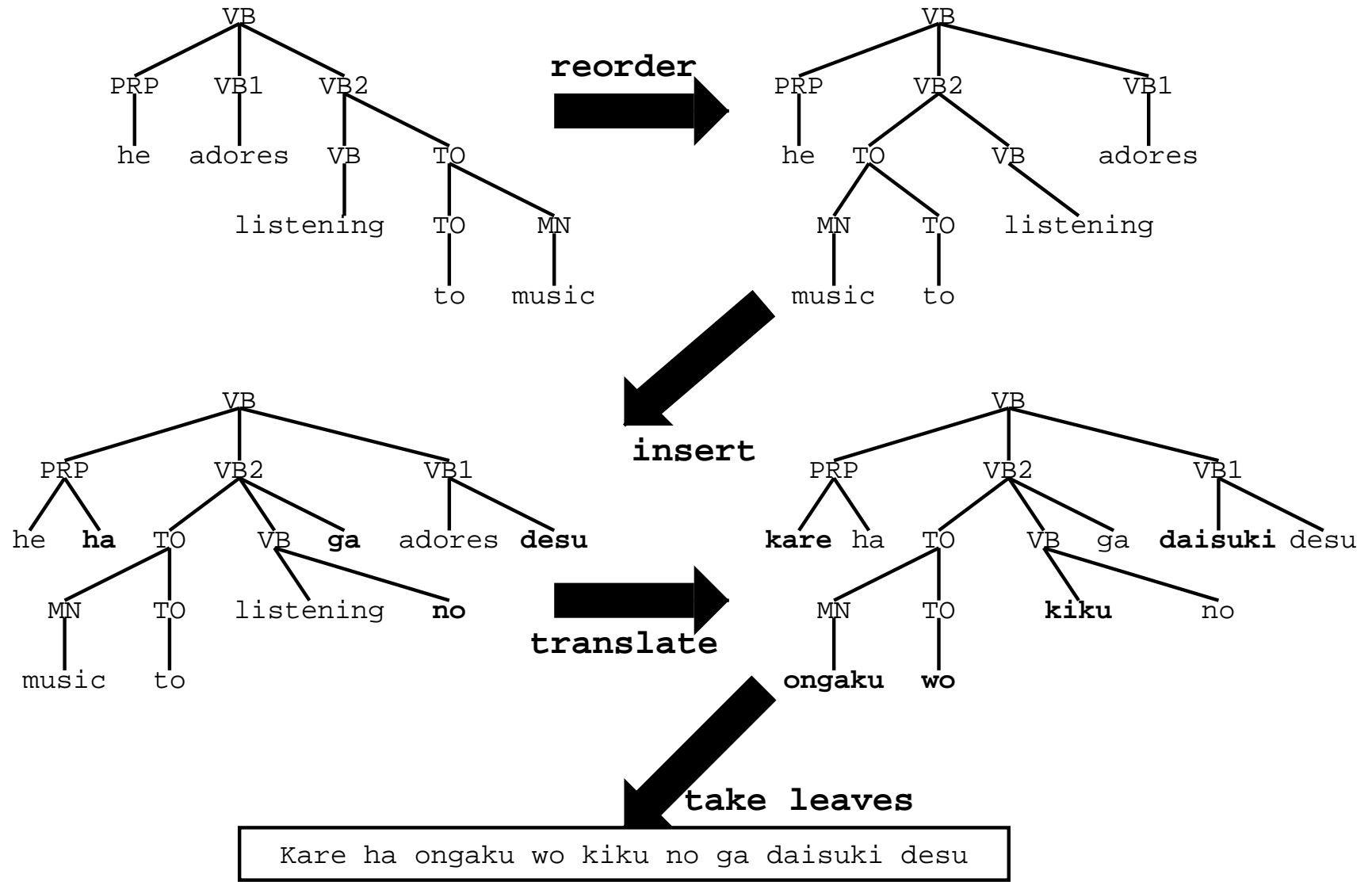
- Do Foreign and English have Same Dependency Structure?
- Direct Correspondence Assumption [Hwa et al., 2002]
 - empirical study (by projection) of Chinese-English parallel corpus
 - even with modifications, only 67% precision/recall
 - more structure could be preserved, if tried

String to Tree Translation



- Use of English Syntax Trees [Yamada and Knight, 2001]
 - exploit rich resources on the English side
 - obtained with statistical parser [Collins, 1997]
 - flattened tree to allow more reorderings
 - works well with syntactic language model

Yamada and Knight [2001]



Crossings

- Do English Trees Match Foreign Strings?
- Crossings between French-English [Fox, 2002]
 - 0.29-6.27 per sentence, depending on how it is measured
- Can be Reduced by
 - flattening tree, as done by [Yamada and Knight, 2001]
 - detecting phrasal translation
 - special treatment for small number of constructions
- Most Coherence between Dependency Structures

Full Syntactic/Semantic Translation

- Existing Systems Hybrid Rule-Based / Statistical
 - Czech-English [Cmejrek et al., 2003]
 - Spanish-English [Habash, 2002]
- Performance Below Phrase-Based Statistical Systems
- Why is it so Hard?
 - loss of good phrasal translations [Koehn et al., 2003]
 - lack of foreign syntactic parsers
 - differences in syntactic structure
 - semantic transfer hard to learn (no parallel data)

Outline

- Data
- Evaluation
- Introduction to Statistical Machine Translation
- Translation Model
- **Language Model**
- Decoding Algorithm
- New Directions: Divide and Conquer
- Available Resources

Language Model

- Goal of the Language Model:
Detect good English

Language Model

- What is Good English?
- Standard Technique: Trigram Model
 - multiplication of trigram probabilities
 - $p(\text{witch}|\text{the green}) > p(\text{green}|\text{the witch})$

Mary did not slap the green witch

Mary => $p(\text{Mary})$

Mary did => $p(\text{did}|\text{Mary})$

Mary did not => $p(\text{not}|\text{Mary did})$

did not slap => $p(\text{slap}|\text{did not})$

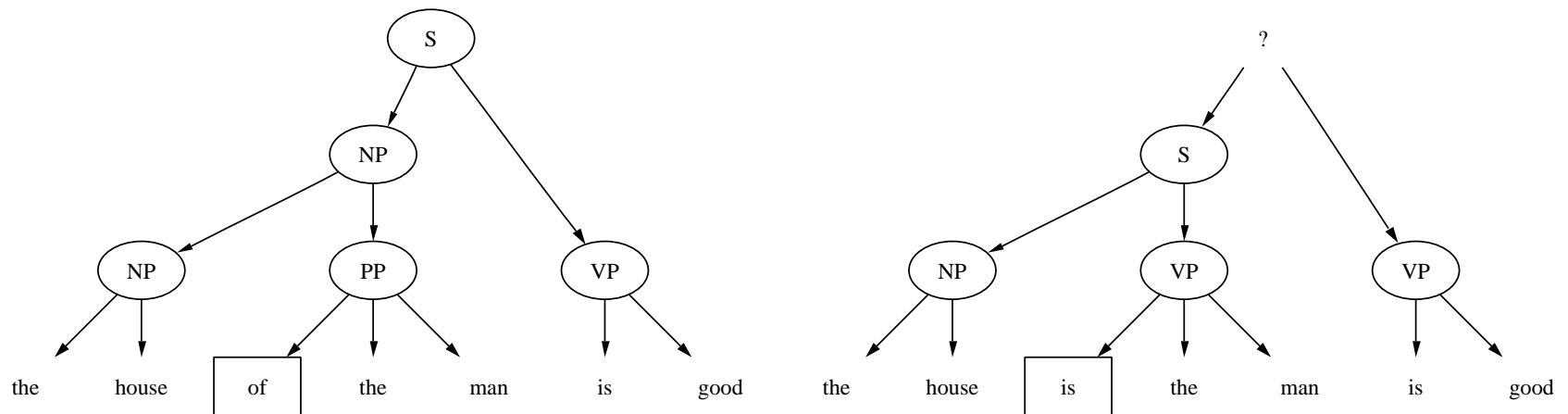
not slap the => $p(\text{the}|\text{not slap})$

slap the green => $p(\text{green}|\text{slap the})$

the green witch => $p(\text{witch}|\text{the green})$

Syntactic Language Model

- Good Syntax Tree \rightarrow Good English
- Allows for Long Distance Constraints



- Left Translation Preferred by Syntactic LM

Using Web n-Grams as LM

- n-Grams Seen on Web:

	Human translation	Machine translation
bigrams	99% seen on web	97%
trigrams	97%	92%
4-grams	85%	80%
5-grams	65%	56%
6-grams	44%	32%
7-grams	30%	14%

- Successfully Used Web n-Grams as Feature
[Koehn and Knight, 2003]

Exploiting Non-Parallel Corpora

- Use Frequencies on the Web [Soricut et al., 2002]
 - She has a lot of nerve. (20 Altavista)
 - It has a lot of nerve. (3 Altavista)
- Build Suffix Trees [Munteanu and Marcu, 2002]
- Learn Bilingual Dictionary Weights
[Koehn and Knight, 2000]

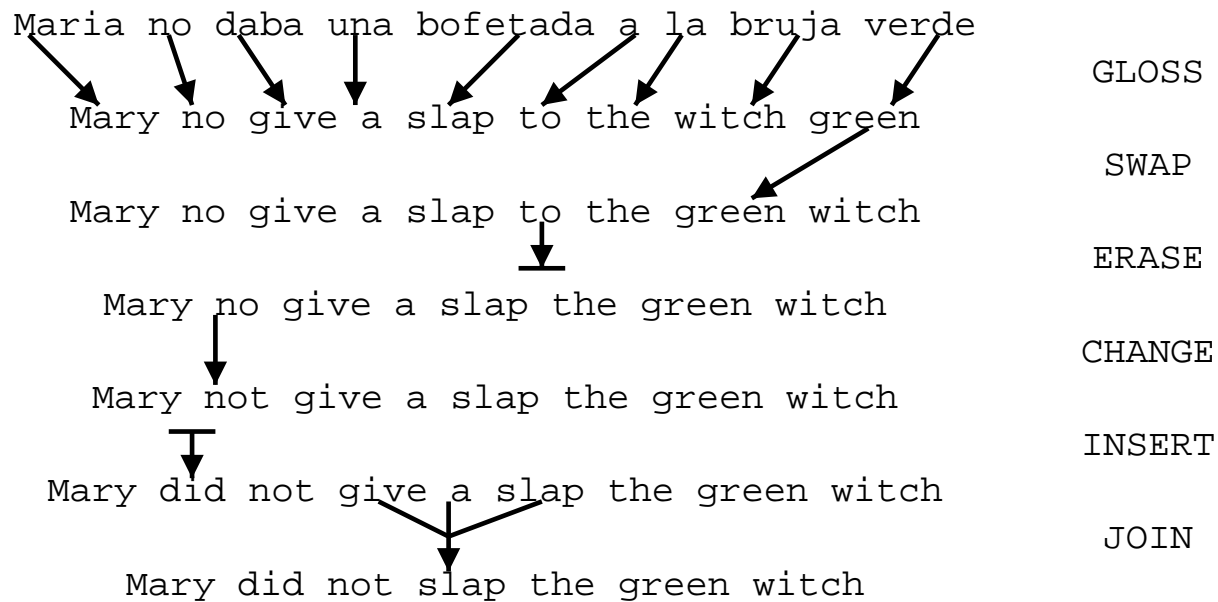
Outline

- Data
- Evaluation
- Introduction to Statistical Machine Translation
- Translation Model
- Language Model
- **Decoding Algorithm**
- New Directions: Divide and Conquer
- Available Resources

Decoding Algorithm

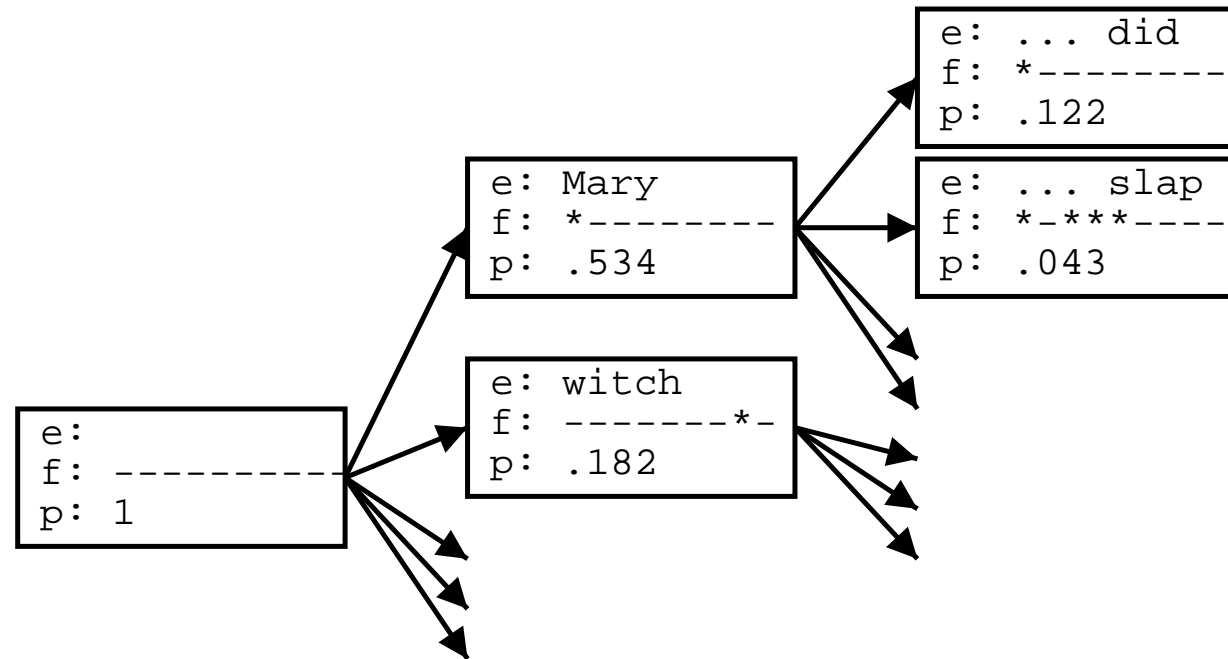
- Goal of the decoding algorithm:
Put models to work, perform the actual translation

Greedy Decoder



- Greedy Hill-climbing [Germann, 2003]
 - start with gloss
 - improve probability with actions
 - use 2-step look-ahead to avoid some local minima

Beam Search Decoding

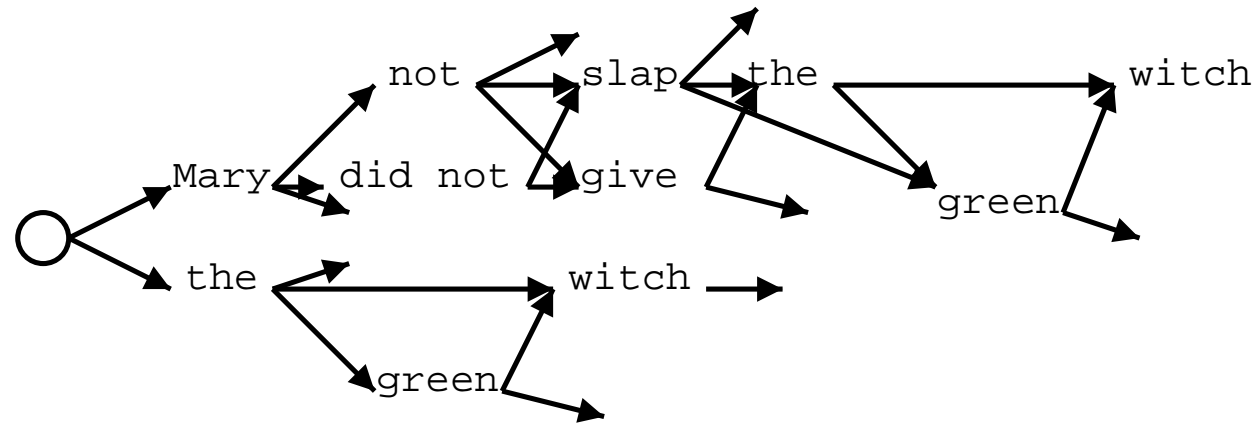


- Build English by Hypothesis Expansion
 - from left to right
 - search space exponential with sentence length
- ⇒ reduction by pruning weak hypothesis

Beam: Search Space Reduction

- Organize Hypotheses into Bins
 - same foreign words covered (still exponential)
 - same number of foreign words covered
 - same number of English words generated
- Prune out Weakest Hypotheses in Each Bin
 - by absolute threshold (keep 100 best)
 - by relative cutoff (only if <0.01 worse than best)
- Future Cost Estimation
 - to have a more realistic comparison of hypothesis
 - compute expected cost of untranslated words
 - add to accumulated cost so far

Beam: Word Graphs



- Word Graphs

- search graph from beam search can be easily converted
- important: hypothesis recombination
- can be mined for n-best lists [Ueffing et al., 2002]

Other Decoding Methods

- Finite State Transducers
 - e.g., [Al-Onaizan and Knight, 1998], [Alshawi et al., 1997]
 - well studied framework, many tools available
- Integer Programming [Germann et al., 2001]
- For String to Tree Model: Parsing
 - see [Yamada and Knight, 2002]
 - uses dynamic programming, similar to chart parsing
 - hypothesis space can be efficiently encoded in forest structure

Outline

- Data
- Evaluation
- Introduction to Statistical Machine Translation
- Translation Model
- Language Model
- Decoding Algorithm
- **New Directions: Divide and Conquer**
- Available Resources

New Directions

- How can we add more knowledge to the process?
 - Define subtasks
 - Maximum entropy framework to include more features

Divide and Conquer

- Named Entities

- names
- numbers
- dates
- quantities

- Noun Phrases

Numbers, Dates, Entities

- Translation Tables for Numbers?

f	e	p(f e)
2003	2003	0.7432
2003	2000	0.0421
2003	year	0.0212
2003	the	0.0175
2003

- Or by Special Handling?

- XML markup of MT input [Germann et al., 2003]
- the revenue for
`<number translate-as='2003'> 2003 </number>`
is higher than ...
- same for dates and quantities
- infinite variety, but simple translation rules

Names

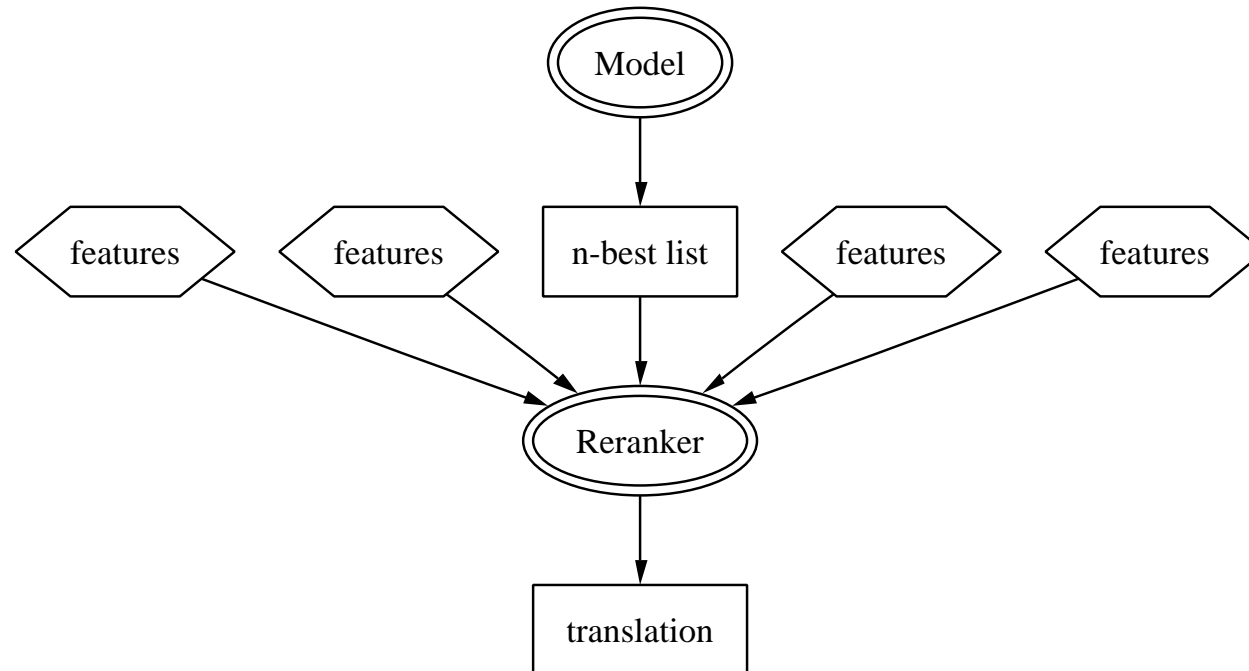
- Often not in Training Corpus
- Require Special Treatment
- Issues
 - recognition of name vs. non-name
 - translation (Defense Department) vs. transliteration (George Bush)
 - especially hard, if different character set (Arabic, Chinese, Cyrillic, ...)
- Phonetic Reasoning and Web Resources

Arabic-English	all	person	organization	location
Sakhr	61%	47%	81%	36%
[Al-Onaizan and Knight, 2002]	73%	64%	87%	51%
Human	75%	68%	95%	42%

Noun Phrases

- Noun Phrases can be Translated in Separation [Koehn and Knight, 2003]
 - German-English: 75% are, 98% can be
 - also other examined languages: Portuguese-E, Chinese-E
 - Definition of NP/PP
 - (informally): maximal phrases that contain at least one noun and no verb
 - (*The permanent tribunal*) is designed to prosecute (*individuals*) (*for genocide, crimes against humanity and other war crimes*) .
 - cover about half of the words, all nouns (largest open word class)
 - shorter, simpler than full sentences
- ⇒ special linguistic modeling, expensive features

Noun Phrases: Re-Ranking



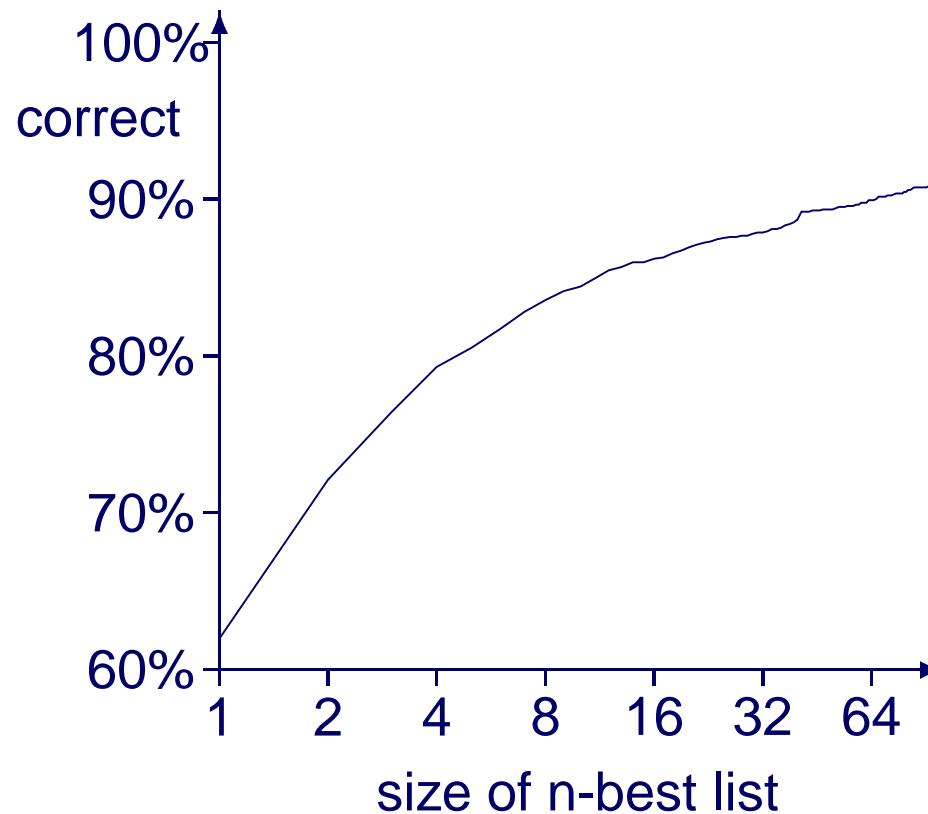
- **Maximum Entropy Reranking**

- allows for variety of features: binary, integer, real-valued
- see also direct maximum entropy models [Och and Ney, 2002]

Noun Phrases: Re-Ranking (2)

- Correct Translations in the n-Best List

→ over 90% accuracy possible with 100-best list reranking



Noun Phrases: Results

- Results for German-English

System	NP/PP Correct	BLEU Full Sentence
IBM Model 4	53.2%	0.172
Phrase Model	58.7%	0.188
Compound Splitting	61.5%	0.195
Re-Estimated Parameters	63.0%	0.197
Web Count Features	64.7%	0.198
Syntactic Features	65.5%	0.199

How Good is Statistical MT?

- Out-of-domain (Sports)

Basketball Network and Valve Promoted More Eastern Second Round

Washington (Afp) new Jersey nets basketball team Thursday again rather than Indian it slipped horseback birds will be Miller of selling your life and hard work, the two extensions to competition after more than 120 109 to Clinton slipped horseback, winning more quarter after the competition for the first round matches of the war, and promoted the second round...

- In-domain (Politics)

The United States and India May Will Be Held in the Past 40 Years the First Joint Military Exercises

(Afp report from new Delhi) India and U. S. will be held in the past 39 years the first joint military exercises in the world's two biggest democracies the cooperative relationship between making milestone.

The Defense Ministry said in a class Indian paratrooper Brigade mid-May and the US Pacific Command of the special units in the well-known far and near the Thai women Maha tomb near joint military exercises.

The two countries will provide air support.

- DARPA Chinese-English task (fairly hard)

- This is actual output of the ISI system