



UPPSALA
UNIVERSITET

Phrase-based Statistical Machine Translation

Sara Stymne

2018-09-20

Some slides from Philipp Koehn



Lecture outline

- What is phrase-based SMT?
- Phrase-based modeling
- Training
- Log-linear models and features
- Tuning (?)
- Factored models (?)



Word-based vs Phrase-based SMT

- Word-based models translate **words** as atomic units
- Phrase-based models translate **phrases** as atomic units



Word-based vs Phrase-based SMT

- Word-based models translate **words** as atomic units
- Phrase-based models translate **phrases** as atomic units
 - A phrase is a continuous sequence of words
 - Not necessarily a linguistic phrase



Advantages of Phrase-based SMT

- Many-to-many translation can handle
 - Non-compositional phrases
kick the bucket – ins Gras beissen (lit: into the grass bite)
 - Compounds
blädderblocksblad – flipchart paper
- Use of local context
 - Local word order
affaires extérieures – external affairs
 - Local agreement issues
det röda blocket – the red block
den röda konen – the red cone



Advantages of Phrase-based SMT II

- Translating phrases helps to reduce translation ambiguities
- Phrases of arbitrary length: sometimes the entire sentence might be covered by a phrase
- Simpler model: no more need to explicitly model the concepts of fertility, insertion and deletion of words



Word-based SMT: Generative Model

Bakom huset hittade polisen en stor mängd narkotika .

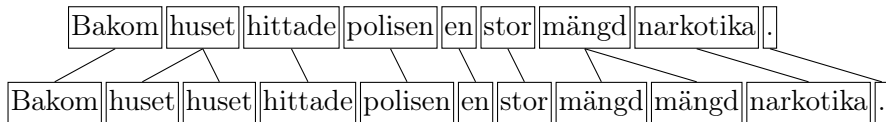


Word-based SMT: Generative Model

Bakom huset hittade polisen en stor mängd narkotika .



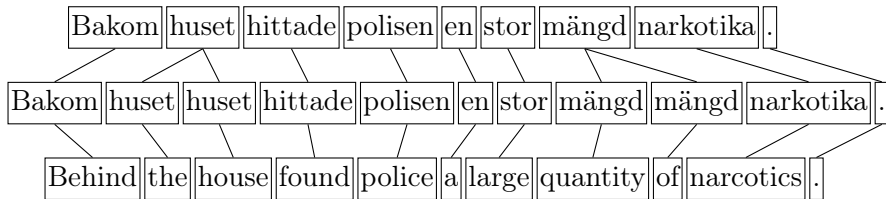
Word-based SMT: Generative Model



1 Fertility



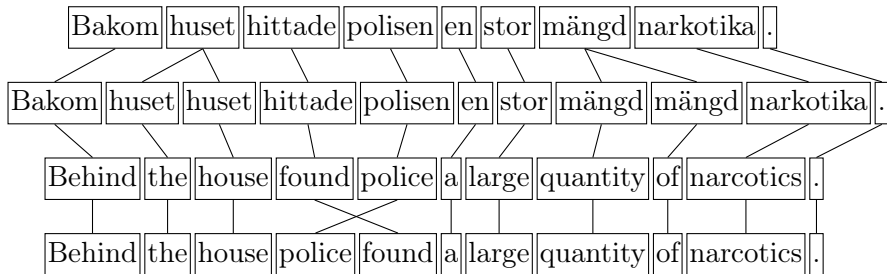
Word-based SMT: Generative Model



- 1 Fertility
- 2 Word translation



Word-based SMT: Generative Model



- 1 Fertility
- 2 Word translation
- 3 Output ordering



Phrase-based SMT: Generative Model

Bakom huset hittade polisen en stor mängd narkotika .



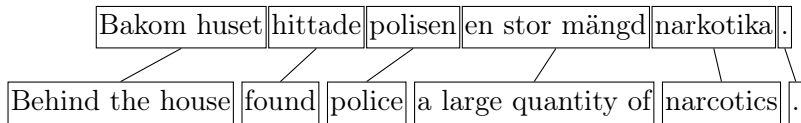
Phrase-based SMT: Generative Model

Bakom huset hittade polisen en stor mängd narkotika.

1 Phrase segmentation



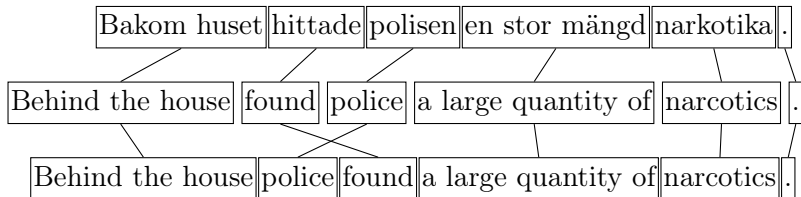
Phrase-based SMT: Generative Model



- 1 Phrase segmentation
- 2 Phrase translation



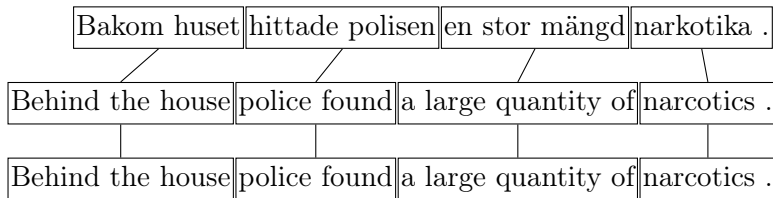
Phrase-based SMT: Generative Model



- 1 Phrase segmentation
- 2 Phrase translation
- 3 Output ordering



Phrase-based SMT: Alternative segmentation



- 1 Phrase segmentation
- 2 Phrase translation
- 3 Output ordering



Phrase translation table

- Main knowledge source: table with phrase translations and their probabilities
- Example: phrase translations for **naturligtvis**

Translation	Probability $\phi(\bar{t} \bar{s})$
of course	0.5
naturally	0.3
of course ,	0.15
, of course ,	0.05



Real example

UPPSALA
UNIVERSITET

- Phrase translations for **begreppet** learned from the Europarl corpus

English	$\phi(\bar{t} \bar{s})$	English	$\phi(\bar{t} \bar{s})$
the	0.226415	the news	0.012816
told	0.169811	the report	0.008544
announcement	0.075472	the information	0.008544
message	0.056604	the back	0.004272
news	0.056604	the suspension	0.004272
information	0.037736	the death	0.004272
informed	0.037736	this announcement	0.002848
learnt	0.037736	this news	0.002136
peace of mind by ensuring	0.027778	a message	0.001539
insight	0.018868	his answer	0.000356
the announcement	0.017088	were told	0.000229
the message	0.012816	the back and	2.917e-05



Real example

UPPSALA
UNIVERSITET

- Phrase translations for **begreppet** learned from the Europarl corpus

English	$\phi(\bar{t} \bar{s})$	English	$\phi(\bar{t} \bar{s})$
the	0.226415	the news	0.012816
told	0.169811	the report	0.008544
announcement	0.075472	the information	0.008544
message	0.056604	the back	0.004272
news	0.056604	the suspension	0.004272
information	0.037736	the death	0.004272
informed	0.037736	this announcement	0.002848
learnt	0.037736	this news	0.002136
peace of mind by ensuring	0.027778	a message	0.001539
insight	0.018868	his answer	0.000356
the announcement	0.017088	were told	0.000229
the message	0.012816	the back and	2.917e-05

- lexical variation (announcement, message, news, told, ...)



Real example

UPPSALA
UNIVERSITET

- Phrase translations for **begreppet** learned from the Europarl corpus

English	$\phi(\bar{t} \bar{s})$	English	$\phi(\bar{t} \bar{s})$
the	0.226415	the news	0.012816
told	0.169811	the report	0.008544
announcement	0.075472	the information	0.008544
message	0.056604	the back	0.004272
news	0.056604	the suspension	0.004272
information	0.037736	the death	0.004272
informed	0.037736	this announcement	0.002848
learnt	0.037736	this news	0.002136
peace of mind by ensuring	0.027778	a message	0.001539
insight	0.018868	his answer	0.000356
the announcement	0.017088	were told	0.000229
the message	0.012816	the back and	2.917e-05

- lexical variation (announcement, message, news, told, ...)
- Morphological variation (information, informed)



Real example

- Phrase translations for **begreppet** learned from the Europarl corpus

English	$\phi(\bar{t} \bar{s})$	English	$\phi(\bar{t} \bar{s})$
the	0.226415	the news	0.012816
told	0.169811	the report	0.008544
announcement	0.075472	the information	0.008544
message	0.056604	the back	0.004272
news	0.056604	the suspension	0.004272
information	0.037736	the death	0.004272
informed	0.037736	this announcement	0.002848
learnt	0.037736	this news	0.002136
peace of mind by ensuring	0.027778	a message	0.001539
insight	0.018868	his answer	0.000356
the announcement	0.017088	were told	0.000229
the message	0.012816	the back and	2.917e-05

- lexical variation (announcement, message, news, told, ...)
- Morphological variation (information, informed)
- Included function words (the, a, were, this)



Real example

UPPSALA
UNIVERSITET

- Phrase translations for **begreppet** learned from the Europarl corpus

English	$\phi(\bar{t} \bar{s})$	English	$\phi(\bar{t} \bar{s})$
the	0.226415	the news	0.012816
told	0.169811	the report	0.008544
announcement	0.075472	the information	0.008544
message	0.056604	the back	0.004272
news	0.056604	the suspension	0.004272
information	0.037736	the death	0.004272
informed	0.037736	this announcement	0.002848
learnt	0.037736	this news	0.002136
peace of mind by ensuring	0.027778	a message	0.001539
insight	0.018868	his answer	0.000356
the announcement	0.017088	were told	0.000229
the message	0.012816	the back and	2.917e-05

- lexical variation (announcement, message, news, told, ...)
- Morphological variation (information, informed)
- Included function words (the, a, were, this)
- Noise (the, the back and)



Phrases

- Model is not limited to linguistic phrases
(noun phrases, verb phrases, prepositional phrases)
- Example of useful non-linguistic phrases:
 - det finns – there is/are
 - put off – skjuta upp



- Model is not limited to linguistic phrases
(noun phrases, verb phrases, prepositional phrases)
- Example of useful non-linguistic phrases:
 - det finns – there is/are
 - put off – skjuta upp
- Experiments have shown that limitation to only linguistic phrases hurts quality



Probabilistic model

- Bayes rule

$$\begin{aligned} t_{best} &= \arg \max_t p(t|s) \\ &= \arg \max_t p(s|t)p_{LM}(t) \end{aligned} \tag{1}$$

- translation model: $p(s|t)$
- language model: $p_{LM}(t)$



Probabilistic model

- Bayes rule

$$\begin{aligned} t_{best} &= \arg \max_t p(t|s) \\ &= \arg \max_t p(s|t)p_{LM}(t) \end{aligned} \tag{1}$$

- translation model: $p(s|t)$

- language model: $p_{LM}(t)$

- Decomposition of translation model

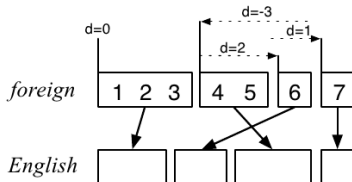
$$p(\bar{s}_1^I | \bar{t}_1^I) = \prod_{i=1}^I \phi(\bar{s}_i | \bar{t}_i) d(\text{start}_i - \text{end}_{i-1} - 1)$$

- phrase translation model: ϕ

- reordering probability: d



Distance-based reordering



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

- distance = $\text{start}_i - \text{end}_{i-1} - 1$
- Scoring function: $d(x) = \alpha^{|x|}$ – exponential with distance



Learning a Phrase Translation Table

- Task: learn the model from a parallel corpus
- Three stages:
 - Word alignment
 - Extraction of phrase pairs
 - Scoring of phrase pairs



Word alignment

UPPSALA
UNIVERSITET

	nyss	hade	jag	precis	tappat	bort	glassen
a	■						
moment	■						
ago	■						
I			■				
had		■					
just				■			
lost					■	■	
my							
ice							■
cream							■



Obtaining a word alignment

- Using some word alignment software and models, often:
 - GIZA++
 - IBM1 – HMM – IBM3 – IBM4



Obtaining a word alignment

- Using some word alignment software and models, often:
 - GIZA++
 - IBM1 – HMM – IBM3 – IBM4
- Such models are directional
 - Gives 1–N links
 - Does not give M–1 or M–N links



Obtaining a word alignment

- Using some word alignment software and models, often:
 - GIZA++
 - IBM1 – HMM – IBM3 – IBM4
- Such models are directional
 - Gives 1–N links
 - Does not give M–1 or M–N links
- We want all types of links!
- Solution: symmetrize directional alignments



Directional word alignment

	nyss	hade	jag	precis	tappat	bort	glassen
a	■						
moment	■						
ago	■						
I			■				
had		■					
just				■			
lost					■		
my							
ice							■
cream							■

■ En-Sv (M-1)



Directional word alignment

	nyss	hade	jag	precis	tappat	bort	glassen
a							
moment	■						
ago							
I			■				
had		■					
just				■			
lost					■	■	
my							
ice							
cream							■


■ En-Sv (M-1)


■ Sv-En (1-N)



Word alignment – symmetrization

	nyss	hade	jag	precis	tappat	bort	glassen
a	En-Sv (M-1)						
moment	Both						
ago	En-Sv (M-1)						
I			Both				
had		Both					
just				Both			
lost					Both	Sv-En (1-N)	
my							
ice							En-Sv (M-1)
cream							Both

 En-Sv (M-1)

 Sv-En (1-N)

 Both



Word alignment – intersection

	nyss	hade	jag	precis	tappat	bort	glassen
a							
moment	■						
ago							
I			■				
had		■					
just				■			
lost					■		
my							
ice							
cream							■

Intersection



Word alignment – union

UPPSALA
UNIVERSITET

	nyss	hade	jag	precis	tappat	bort	glassen
a	■						
moment	■						
ago	■						
I			■				
had		■					
just				■			
lost					■	■	
my							
ice							■
cream							■

Union



Word alignment symmetrization

- Intersection: too few links
- Union: too many links

English–Swedish alignment (Holmqvist, 2008)

	Precision	Recall
Intersection	90	75
Union	60	91



Word alignment symmetrization 2

- Solution: use heuristics!



Word alignment symmetrization 2

- Solution: use heuristics!
- grow-diag-final-and
 - Start with the intersection
 - Add links from the union:
 - Add diagonally adjacent links
 - Add links for unaligned words in a final step



Word alignment symmetrization 2

- Solution: use heuristics!
- grow-diag-final-and
 - Start with the intersection
 - Add links from the union:
 - Add diagonally adjacent links
 - Add links for unaligned words in a final step

	Precision	Recall
Intersection	90	75
Union	60	91
grow-diag-final-and	70	88



Word alignment symmetrization 2

- Solution: use heuristics!
- grow-diag-final-and
 - Start with the intersection
 - Add links from the union:
 - Add diagonally adjacent links
 - Add links for unaligned words in a final step
- There are other symmetrization heuristics and word alignment methods
- In general: high recall is important for good SMT quality (with large training data)

	Precision	Recall
Intersection	90	75
Union	60	91
grow-diag-final-and	70	88



Extracting phrase pairs

- Extract phrase pairs that are consistent with word alignments

	Nyss	hade	jag	precis	tappat	bort	glassen
A	■						
moment	■						
ago	■						
I			■				
had		■					
just				■			
lost					■	■	
my					■	■	
ice							■
cream							■



Extracting phrase pairs

- Extract phrase pairs that are consistent with word alignments

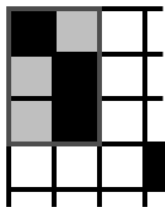
	Nyss	hade	jag	precis	tappat	bort	glassen
A	■						
moment	■						
ago	■						
I			■				
had		■					
just			■				
lost				■			
my				■	■		
ice					■		
cream							■

just lost-precis tappat bort



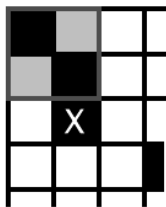
Consistent phrase pairs

All words of the phrase pairs have to align to each other



consistent

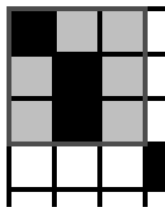
ok



inconsistent

violated

one alignment
point outside



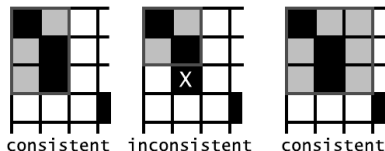
consistent

ok

unaligned
word is fine



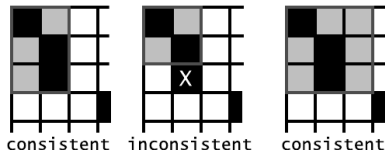
Phrase extraction definition



A phrase pair (\bar{t}, \bar{s}) is consistent with an alignment A , if all words s_1, \dots, s_m in \bar{s} that have alignment points in A have these with words t_1, \dots, t_n in \bar{t} and vice versa and at least one word in \bar{t} is aligned to at least one word in \bar{s}



Phrase extraction definition



A phrase pair (\bar{t}, \bar{s}) is consistent with an alignment A , if all words s_1, \dots, s_m in \bar{s} that have alignment points in A have these with words t_1, \dots, t_n in \bar{t} and vice versa and at least one word in \bar{t} is aligned to at least one word in \bar{s}

$$\begin{aligned}(\bar{t}, \bar{s}) \text{ consistent with } A &\Leftrightarrow \forall t_i \in \bar{t} : (t_i, s_j) \in A \rightarrow s_j \in \bar{s} \\ &\text{AND } \forall s_i \in \bar{s} : (t_i, s_j) \in A \rightarrow t_j \in \bar{t} \\ &\text{AND } \exists t_i \in \bar{t}, s_i \in \bar{s} : (t_i, s_j) \in A\end{aligned}$$



Phrases extracted

UPPSALA
UNIVERSITET

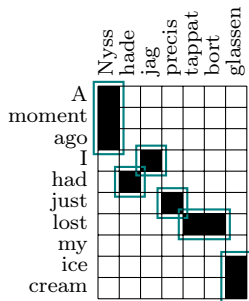
	Nyss	hade	jag	precis	tappat	bort	glassen
A	■						
moment	■						
ago	■						
I			■				
had		■					
just				■			
lost					■	■	
my							
ice							■
cream							■



Phrases extracted

UPPSALA
UNIVERSITET

a moment ago–nyss, I–jag, had–hade, just–precis
lost–tappat bort, ice cream–glassen

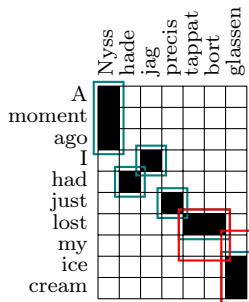




Phrases extracted

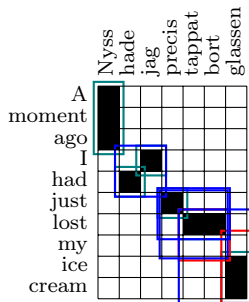
a moment ago–nyss, I–jag, had–hade, just–precis
lost–tappat bort, ice cream–glassen

lost my–tappat bort, my ice cream–glassen





Phrases extracted



a moment ago–nyss, I–jag, had–hade, just–precis
lost–tappat bort, ice cream–glassen

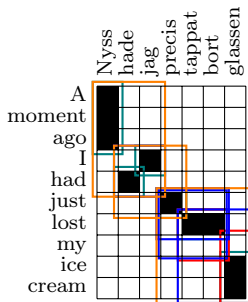
lost my–tappat bort, my ice cream–glassen

I had–hade jag, lost my ice cream–tappat bort glassen
just lost–precis tappat bort, just lost my–precis tappat bort



Phrases extracted

UPPSALA
UNIVERSITET



a moment ago–nyss, I–jag, had–hade, just–precis
lost–tappat bort, ice cream–glassen

lost my–tappat bort, my ice cream–glassen

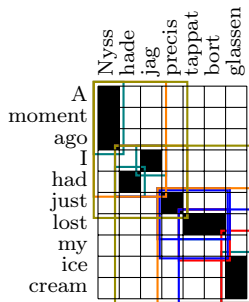
I had–hade jag, lost my ice cream–tappat bort glassen
just lost–precis tappat bort, just lost my–precis tappat bort

a moment ago I had–nyss hade jag, I had just–hade jag precis
just lost my ice cream–precis tappat bort glassen



Phrases extracted

UPPSALA
UNIVERSITET



a moment ago–nyss, I–jag, had–hade, just–precis
lost–tappat bort, ice cream–glassen

lost my–tappat bort, my ice cream–glassen

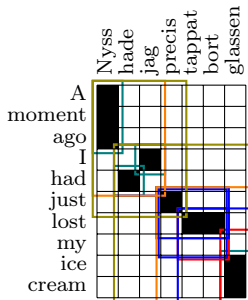
I had–hade jag, lost my ice cream–tappat bort glassen
just lost–precis tappat bort, just lost my–precis tappat bort

a moment ago I had–nyss hade jag, I had just–hade jag precis
just lost my ice cream–precis tappat bort glassen

a moment ago I had just–nyss hade jag precis
I had just lost my ice cream–hade jag precis tappat bort glassen



Phrases extracted



a moment ago–nyss, I–jag, had–hade, just–precis
lost–tappat bort, ice cream–glassen

lost my–tappat bort, my ice cream–glassen

I had–hade jag, lost my ice cream–tappat bort glassen
just lost–precis tappat bort, just lost my–precis tappat bort

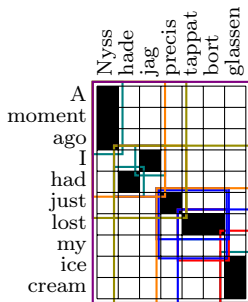
a moment ago I had–nyss hade jag, I had just–hade jag precis
just lost my ice cream–precis tappat bort glassen

a moment ago I had just–nyss hade jag precis
I had just lost my ice cream–hade jag precis tappat bort glassen

...



Phrases extracted



a moment ago–nyss, I–jag, had–hade, just–precis
lost–tappat bort, ice cream–glassen

lost my–tappat bort, my ice cream–glassen

I had–hade jag, lost my ice cream–tappat bort glassen
just lost–precis tappat bort, just lost my–precis tappat bort

a moment ago I had–nyss hade jag, I had just–hade jag precis
just lost my ice cream–precis tappat bort glassen

a moment ago I had just–nyss hade jag precis
I had just lost my ice cream–hade jag precis tappat bort glassen

...

a moment ago I had just lost my ice cream–
nyss hade jag precis tappat bort glassen



Phrase extraction exercise

Exercise on phrase pair extraction!



Scoring phrase translations

- Phrase pair extraction: collect all phrase pairs from the data
- Phrase pair scoring: assign probabilities to phrase translations



Scoring phrase translations

- Phrase pair extraction: collect all phrase pairs from the data
- Phrase pair scoring: assign probabilities to phrase translations
- Score by relative frequency:

$$\phi(\bar{t}|\bar{s}) = \frac{\text{count}(\bar{s}, \bar{t})}{\sum_{\bar{t}_i} \text{count}(\bar{s}, \bar{t}_i)}$$



Scoring phrase translations

- Phrase pair extraction: collect all phrase pairs from the data
- Phrase pair scoring: assign probabilities to phrase translations
- Score by relative frequency:

$$\phi(\bar{t}|\bar{s}) = \frac{\text{count}(\bar{s}, \bar{t})}{\sum_{\bar{t}_i} \text{count}(\bar{s}, \bar{t}_i)}$$

- Potentially improve scoring by smoothing



Size of the phrase table

- Phrase translation table typically much bigger than corpus
- Limit the length of phrase pairs (often to 7 tokens)



Size of the phrase table

- Phrase translation table typically much bigger than corpus
- Limit the length of phrase pairs (often to 7 tokens)
- Too big to store in memory?
 - Store on disk
 - Use smart data structures



Size of the phrase table

- Phrase translation table typically much bigger than corpus
- Limit the length of phrase pairs (often to 7 tokens)
- Too big to store in memory?
 - Store on disk
 - Use smart data structures
- Prune phrase table – i.e., remove non-useful phrase pairs
 - Limit translation options for each phrase (often to 20–30)
 - Prune table based on statistics, such as χ^2



Size of the phrase table

- Phrase translation table typically much bigger than corpus
- Limit the length of phrase pairs (often to 7 tokens)
- Too big to store in memory?
 - Store on disk
 - Use smart data structures
- Prune phrase table – i.e., remove non-useful phrase pairs
 - Limit translation options for each phrase (often to 20–30)
 - Prune table based on statistics, such as χ^2
- Filter phrase table towards a test set



Weighted models

- Described model consists of three sub-models:
 - Phrase translation models $\phi(\bar{s}|\bar{t})$
 - Reordering model d
 - Language model $p_{LM}(t)$

$$t_{best} = \arg \max_t \prod_{i=1}^I \phi(\bar{s}_i|\bar{t}_i) d(\text{start}_i - \text{end}_{i-1} - 1) \prod_{i=1}^{|t|} p_{LM}(t_i|t_{i-(n-1)} \dots t_{i-1})$$



Weighted models

- Described model consists of three sub-models:
 - Phrase translation models $\phi(\bar{s}|\bar{t})$
 - Reordering model d
 - Language model $p_{LM}(t)$

$$t_{best} = \arg \max_t \prod_{i=1}^I \phi(\bar{s}_i|\bar{t}_i) d(\text{start}_i - \text{end}_{i-1} - 1) \prod_{i=1}^{|t|} p_{LM}(t_i|t_{i-(n-1)} \dots t_{i-1})$$

- Some sub-models may be more important than others



Weighted models

- Described model consists of three sub-models:
 - Phrase translation models $\phi(\bar{s}|\bar{t})$
 - Reordering model d
 - Language model $p_{LM}(t)$

$$t_{best} = \arg \max_t \prod_{i=1}^I \phi(\bar{s}_i|\bar{t}_i) d(\text{start}_i - \text{end}_{i-1} - 1) \prod_{i=1}^{|t|} p_{LM}(t_i|t_{i-(n-1)} \dots t_{i-1})$$

- Some sub-models may be more important than others
- Add weights $\lambda_\phi, \lambda_d, \lambda_{LM}$



Weighted models

- Described model consists of three sub-models:
 - Phrase translation models $\phi(\bar{s}|\bar{t})$
 - Reordering model d
 - Language model $p_{LM}(t)$

$$t_{best} = \arg \max_t \prod_{i=1}^I \phi(\bar{s}_i|\bar{t}_i)^{\lambda_\phi} d(\text{start}_i - \text{end}_{i-1} - 1)^{\lambda_d} \prod_{i=1}^{|t|} p_{LM}(t_i|t_{i-(n-1)} \dots t_{i-1})^{\lambda_{LM}}$$

- Some sub-models may be more important than others
- Add weights $\lambda_\phi, \lambda_d, \lambda_{LM}$



Log-linear models

- Such a weighted model is a log-linear model:

$$p(x) = \exp \sum_{i=1}^n \lambda_i h_i(x)$$

- Our feature functions:
 - three feature functions $n = 3$
 - random variable $x = (s, t, \text{start}, \text{end})$
 - feature function $h_1 = \log \phi$
 - feature function $h_2 = \log d$
 - feature function $h_3 = \log p_{LM}$



Weighted model as a log-linear model

$$p(t, a|s) = \exp(\lambda_\phi \sum_{i=1}^I \log \phi(\bar{s}_i | \bar{t}_i) + \\ \lambda_d \sum_{i=1}^I \log d(\text{start}_i - \text{end}_{i-1} - 1) + \\ \lambda_{LM} \sum_{i=1}^{|t|} \log p_{LM}(t_i | t_{i-(n-1)} \dots t_{i-1}))$$



More feature functions

$$t^* = \arg \max_t \sum_i \lambda_i h_i(s, t)$$

- Easy and useful to add more feature functions



More feature functions

$$t^* = \arg \max_t \sum_i \lambda_i h_i(s, t)$$

- Easy and useful to add more feature functions
 - Bidirectional alignment probabilities $\phi(\bar{s}|\bar{t})$ and $\phi(\bar{t}|\bar{s})$



More feature functions

$$t^* = \arg \max_t \sum_i \lambda_i h_i(s, t)$$

- Easy and useful to add more feature functions
 - Bidirectional alignment probabilities $\phi(\bar{s}|\bar{t})$ and $\phi(\bar{t}|\bar{s})$
 - Lexical weighting of phrase pairs:



More feature functions

$$t^* = \arg \max_t \sum_i \lambda_i h_i(s, t)$$

- Easy and useful to add more feature functions
 - Bidirectional alignment probabilities $\phi(\bar{s}|\bar{t})$ and $\phi(\bar{t}|\bar{s})$
 - Lexical weighting of phrase pairs:

$$\text{lex}(\bar{t}|\bar{s}, a) = \prod_{i=1}^{\text{length}(\bar{t})} \frac{1}{|\{j|(i, j) \in a\}|} \sum_{\forall (i, j) \in a} w(t_i|s_j)$$



More feature functions

$$t^* = \arg \max_t \sum_i \lambda_i h_i(s, t)$$

- Easy and useful to add more feature functions
 - Bidirectional alignment probabilities $\phi(\bar{s}|\bar{t})$ and $\phi(\bar{t}|\bar{s})$
 - Lexical weighting of phrase pairs:

$$\text{lex}(\bar{t}|\bar{s}, a) = \prod_{i=1}^{\text{length}(\bar{t})} \frac{1}{|\{j|(i, j) \in a\}|} \sum_{\forall (i, j) \in a} w(t_i|s_j)$$

- Useful since rare phrase pairs have unreliable probability estimates



More feature functions 2

- Language model has a bias towards short translations
 - word count: $wc(t) = \log |t|^\omega$



More feature functions 2

- Language model has a bias towards short translations
 - word count: $wc(t) = \log |t|^\omega$
- We may prefer finer or coarser segmentations
 - phrase count: $pc(t) = \log |I|^\rho$



More feature functions 2

- Language model has a bias towards short translations
 - word count: $wc(t) = \log |t|^\omega$
- We may prefer finer or coarser segmentations
 - phrase count: $pc(t) = \log |I|^\rho$
- Multiple language models



More feature functions 2

- Language model has a bias towards short translations
 - word count: $wc(t) = \log |t|^\omega$
- We may prefer finer or coarser segmentations
 - phrase count: $pc(t) = \log |I|^\rho$
- Multiple language models
- Other knowledge sources



More feature functions 2

- Language model has a bias towards short translations
 - word count: $wc(t) = \log |t|^\omega$
- We may prefer finer or coarser segmentations
 - phrase count: $pc(t) = \log |I|^\rho$
- Multiple language models
- Other knowledge sources
- Lexicalized reordering models



Translation model learning

- Phrase-based model, but trained using word alignments



Translation model learning

- Phrase-based model, but trained using word alignments
- Can we not train a phrase-based model directly?



Translation model learning

- Phrase-based model, but trained using word alignments
- Can we not train a phrase-based model directly?
- Yes, with EM!
 - But, difficult, and often overfits



Translation model learning

- Phrase-based model, but trained using word alignments
- Can we not train a phrase-based model directly?
- Yes, with EM!
 - But, difficult, and often overfits
- The pipeline I have shown is the state-of-the-art



- SMT toolkit
- Free, open source
- Implements several models:
 - Phrase-based
 - Hierarchical
 - Syntax-based
- Decoding (next lecture)
- Training pipeline
 - Training translation models
 - Training language models
 - Optimizing feature weights



Moses training pipeline

- 1 Prepare data
- 2 Run GIZA to create one-way alignments
- 3 Symmetrize alignment
- 4 Calculate lexical translation probabilities
- 5 Extract phrases
- 6 Score phrases
- 7 Train reordering model
- 8 (Train generation model)
- 9 Create configuration file



Moses training pipeline

- Preprocess data (tokenization, casing, et.c.)

- 1 Prepare data
- 2 Run GIZA to create one-way alignments
- 3 Symmetrize alignment
- 4 Calculate lexical translation probabilities
- 5 Extract phrases
- 6 Score phrases
- 7 Train reordering model
- 8 (Train generation model)
- 9 Create configuration file

- Tune the system



Tuning – Optimizing feature weights

- How do we learn the best weights λ_i ?



Tuning – Optimizing feature weights

- How do we learn the best weights λ_i ?
- Optimize the weights on a small corpus



Tuning – Optimizing feature weights

- How do we learn the best weights λ_i ?
- Optimize the weights on a small corpus
- Called **Tuning**



Coming up

- Next week:
 - Monday September 24: deadline for wishing for project topics
 - Monday: Lecture on decoding
 - Tuesday: Assignment 4: PBSMT and Moses
 - Wednesday: Lab 1 (not 5LN718)
 - Thursday: Lecutre cancelled
- Later:
 - Assignment 5: decoding and document-level MT
 - Start working on your projects
 - Lectures and assignment on NMT
 - Guest lecture