## Practical introduction to Statistical Machine Translation

Cristina España i Bonet

GF Summer School, Fraueninsel

23rd August 2013







#### Part I: SMT background

 $\sim$  120min







#### Part II: SMT experiments

from 30min to ...



8 Evaluation system

Part III: MT evaluation

 $\sim$  45min

## Part I

# SMT background

### 1 Introduction

#### 2 Basics

3 Components

- 4 The log-linear model
- 5 Beyond standard SMT











#### Introduction Empirical Machine Translation

Empirical MT relies on aligned corpora



Introduction Empirical Machine Translation

#### Empirical MT relies on aligned corpora





#### Empirical MT relies on large parallel aligned corpora



Som a l'any 50 abans de Crist. Tota la Gàl-lia és ocupada pels romans... Tota? No! Un llogaret del Nord habitat per gals indomables rebutja una i altra vegada ferotgement l'invassor. La vida doncs no és gens planera per als legionaris romans dels petits campaments de Babaòrum, Aquàrium, Laundànum i Petibònum...

The year is 50 B.C. Gaul is entirely occupied by the Romans. Well, not entirely... One small village of indomitable Gauls still holds out against the invaders. And life is not easy for the Roman legionaries who garrison the fortified camps of Totorum, Aquarium, Laudanum and Compendium...

#### Empirical MT relies on large parallel aligned corpora



Som a l'any 50 abans de Crist. Tota la Gàl·lia és ocupada pels romans... Tota? No! Un llogaret del Nord habitat per gals indomables rebutja una i altra vegada ferotgement l'invassor. La vida doncs no és gens planera per als legionaris romans dels petits campaments de Babaòrum, Aquàrium, Laundànum i Petibònum...

The year is 50 B.C. Gaul is entirely occupied by the Romans. Well, not entirely... One small village of indomitable Gauls still holds out against the invaders. And life is not easy for the Roman legionaries who garrison the fortified camps of Totorum, Aquarium, Laudanum and Compendium...

#### Empirical MT relies on large parallel aligned corpora

Som a l'any 50 abans de Crist. Tota la Gàl·lia és ocupada pels romans... Tota? No! Un llogaret del Nord habitat per gals indomables rebutja una i altra vegada ferotgement l'invassor. La vida doncs no és gens planera per als legionaris romans dels petits campaments de Babaòrum, Aquàrium, Laundànum i Petibònum...

Astèrix. És l'heroic petit guerrer d'aquestes aventures, viu com una centella i enginyosament astut. Per això sempre li són encomanades les missions més perilloses. Extrau la seva terrorífica força de la beguda màgica inventada pel druida Panoràmix.

Obèlix. És l'antic inseparable d'Astèrix. Fa de repartidor de menhirs i li agrada d'allò més la carn de porc senglar. És capaç d'abandonar-ho tot per tal de seguir Astèrix en una nova aventura. Sobretot si no hi manquen els senglars i fortes batusses.

Copdegarròtix. És el cap de la tribu. Majestuós, valent i desconfiat alhora, el vell guerre ès respectat pels seus homes i temut pels seus enemics. Tan sols una cosa li fa por: que el cel li pugui caure damunt del cap! Pero, tal com ell mateix acostuma a dir, "Qui dia passa, any empeny!". The year is 50 B.C. Gaul is entirely occupied by the Romans. Well, not entirely... One small village of indomitable Gauls still holds out against the invaders. And life is not easy for the Roman legionaries who garrison the fortified camps of Totorum, Aquarium, Laudanum and Compendium...

Asterix, the hero of these adventures. A shrewd, cunning little warrior; all perilous missions are immediately entrusted to him. Asterix gets his superhuman strength from the magic potion brewed by the druid Getafix...

Obelix, Asterix's inseparable friend. A menhir delivery-man by trade; addicted to wild boar. Obelix is always ready to drop everything and go off on a new adventure with Asterix - so long as there's wild boar to eat, and plenty of fighting.

Finally, Vitalstitistix, the chief of the tribe. Majestic, brave and hot-tempered, the old warrior is respected by his men and feared by his enemies. Vitalstitistix himself has only one fear, he is afraid the sky may fall on his head tomorrow. But as he always says, "Tomorrow never comes". -

#### Aligned parallel corpora: Numbers

#### Corpora

Corpus	# segments (app.)	# words (app.)
JRC-Acquis	$1.0\cdot 10^6$	$30\cdot 10^6$
Europarl	$2.0\cdot 10^6$	$55\cdot 10^6$
United Nations	$10.7\cdot 10^{6}$	$300\cdot 10^6$

#### **Books**

Title	# words (approx.)
The Bible	$0.8 \cdot 10^{6}$
The Dark Tower series	$1.2 \cdot 10^{6}$
Encyclopaedia Britannica	44 · 10 <sup>6</sup>

-

#### Aligned parallel corpora: Numbers

#### Corpora

Corpus	# segments (app.)	# words (app.)
JRC-Acquis	$1.0\cdot 10^6$	$30\cdot 10^6$
Europarl	$2.0\cdot 10^6$	$55\cdot 10^6$
United Nations	$10.7\cdot 10^6$	$300\cdot 10^6$

#### Books

Title	<pre># words (approx.)</pre>
The Bible	$0.8\cdot 10^6$
The Dark Tower series	$1.2\cdot 10^6$
Encyclopaedia Britannica	$44 \cdot 10^6$

#### 

#### WMT13 parallel data

Corpus	# segments	# tokens	
Europarl ENG	1,928,274	52,048,855	
Europarl SPA	1,928,274	53,996,661	
News Commentary ENG	155,615	3,901,839	
News Commentary SPA	155,615	4,364,802	
United Nations ENG	10,749,388	283,672,192	
United Nations SPA	10,749,388	318,045,340	
Total (ENG+SPA)	25,666,554	716,029,689	

http://www.statmt.org/wmt13/translation-task.html



#### In practice

Shows real examples of the previous theory, always from freely available data/software:

- Data: www.statmt.org/wmt13/
- Software: SRILM, GIZA++ & Moses

Standard tools, but not exclusive

More on the hands on!

### 1 Introduction



3 Components

- 4 The log-linear model
- 5 Beyond standard SMT











#### The Noisy Channel as a statistical approach to translation:





#### The Noisy Channel as a statistical approach to translation:





#### The Noisy Channel as a statistical approach to translation:



#### SMT, basics The Noisy Channel approach



Mathematically:

P(e|f)

#### SMT, basics The Noisy Channel approach



Mathematically:

$$P(e|f) = \frac{P(e) P(f|e)}{P(f)}$$

 $T(f) = \hat{e} = \operatorname{argmax}_{e} P(e|f) = \operatorname{argmax}_{e} P(e) P(f|e)$ 



$$T(f) = \hat{e} = \operatorname{argmax}_{e} P(e) P(f|e)$$

#### Language Model

- Takes care of fluency in the target language
- Data: corpora in the target language

#### Translation Model

- Lexical correspondence between languages
- Data: aligned corpora in source and target languages

argmax

• Search done by the *decoder* 

$$T(f) = \hat{e} = \operatorname{argmax}_{e} P(e) P(f|e)$$

Language Model

- Takes care of fluency in the target language
- Data: corpora in the target language

#### Translation Model

- Lexical correspondence between languages
- Data: aligned corpora in source and target languages

argmax

• Search done by the *decoder* 

$$T(f) = \hat{e} = \operatorname{argmax}_{e} P(e) P(f|e)$$

Language Model

- Takes care of fluency in the target language
- Data: corpora in the target language

Translation Model

- Lexical correspondence between languages
- Data: aligned corpora in source and target languages

argmax

• Search done by the *decoder* 

#### Introduction





- Language model
- Translation model
- Decoder
- The log-linear model
- 5 Beyond standard SMT

#### SMT, components The language model P(e)

Language model

$$T(f) = \hat{e} = \operatorname{argmax}_{e} \frac{P(e) P(f|e)}{P(f|e)}$$
  
Estimation of how probable a sentence is.

Naïve estimation on a corpus with N sentences:

Frequentist probability of a sentence *e*:

$$P(e) = \frac{N_e}{N_{sentences}}$$

Problem:

Long chains are difficult to observe in corpora.
⇒ Long sentences may have zero probability!

#### SMT, components The language model P(e)

#### Language model

$$T(f) = \hat{e} = \operatorname{argmax}_{e} \frac{P(e) P(f|e)}{P(f|e)}$$
  
Estimation of how probable a sentence is

Naïve estimation on a corpus with N sentences:

Frequentist probability of a sentence e: P(e)

$$P(e) = rac{N_e}{N_{sentences}}$$

Problem:

Long chains are difficult to observe in corpora.
⇒ Long sentences may have zero probability!

#### SMT, components The language model P(e)

Language model

$$T(f) = \hat{e} = \operatorname{argmax}_{e} \frac{P(e) P(f|e)}{P(f|e)}$$
  
Estimation of how probable a sentence is

Naïve estimation on a corpus with N sentences:

Frequentist probability of a sentence e:  $P(e) = \frac{N_e}{N_{sentences}}$ 

Problem:

Long chains are difficult to observe in corpora.
⇒ Long sentences may have zero probability!
#### The n-gram approach

The language model assigns a probability P(e)to a sequence of words  $e \Rightarrow \{w_1, \dots, w_m\}$ .  $P(w_1, \dots, w_m) = \prod_{i=1}^m P(w_i | w_{i-(n-1)}, \dots, w_{i-1})$ 

- The probability of a sentence is the product of the conditional probabilities of each word *w<sub>i</sub>* given the previous ones.
- Independence assumption: the probability of *w<sub>i</sub>* is only conditioned by the *n* previous words.

#### Example, a 4-gram model

#### e: All work and no play makes Jack a dull boy

$$\begin{split} P(e) &= P(\text{All}|\phi, \phi, \phi) \; P(\text{work}|\phi, \phi, \text{All}) \; P(\text{and}|\phi, \text{All}, \text{work}) \\ &= P(\text{no}|\text{All}, \text{work}, \text{and}) \; P(\text{play}|\text{work}, \text{and}, \text{no}) \\ &= P(\text{makes}|\text{and}, \text{no}, \text{play}) P(\text{Jack}|\text{no}, \text{play}, \text{makes}) \\ &= P(\text{a}|\text{play}, \text{makes}, \text{Jack}) P(\text{dull}|\text{makes}, \text{Jack}, \text{a}) \\ &= P(\text{boy}|\text{Jack}, \text{a}, \text{dull}) \end{split}$$

#### Example, a 4-gram model

#### e: All work and no play makes Jack a dull boy

$$\begin{split} P(e) &= P(\texttt{All}|\phi, \phi, \phi) \ P(\texttt{work}|\phi, \phi, \texttt{All}) \ P(\texttt{and}|\phi, \texttt{All}, \texttt{work}) \\ P(\texttt{no}|\texttt{All}, \texttt{work}, \texttt{and}) \ P(\texttt{play}|\texttt{work}, \texttt{and}, \texttt{no}) \\ P(\texttt{makes}|\texttt{and}, \texttt{no}, \texttt{play}) P(\texttt{Jack}|\texttt{no}, \texttt{play}, \texttt{makes}) \\ P(\texttt{a}|\texttt{play}, \texttt{makes}, \texttt{Jack}) P(\texttt{dull}|\texttt{makes}, \texttt{Jack}, \texttt{a}) \\ P(\texttt{boy}|\texttt{Jack}, \texttt{a}, \texttt{dull}) \end{split}$$

#### Example, a 4-gram model

e: All work and no play makes Jack a dull boy

$$\begin{split} P(e) &= P(\texttt{All}|\phi, \phi, \phi) \; P(\texttt{work}|\phi, \phi, \texttt{All}) \; P(\texttt{and}|\phi, \texttt{All}, \texttt{work}) \\ &\quad P(\texttt{no}|\texttt{All}, \texttt{work}, \texttt{and}) \; P(\texttt{play}|\texttt{work}, \texttt{and}, \texttt{no}) \\ &\quad P(\texttt{makes}|\texttt{and}, \texttt{no}, \texttt{play}) P(\texttt{Jack}|\texttt{no}, \texttt{play}, \texttt{makes}) \\ &\quad P(\texttt{a}|\texttt{play}, \texttt{makes}, \texttt{Jack}) P(\texttt{dull}|\texttt{makes}, \texttt{Jack}, \texttt{a}) \\ &\quad P(\texttt{boy}|\texttt{Jack}, \texttt{a}, \texttt{dull}) \end{split}$$

#### Example, a 4-gram model

e: All work and no play makes Jack a dull boy

 $P(e) = P(\text{All}|\phi, \phi, \phi) P(\text{work}|\phi, \phi, \text{All}) P(\text{and}|\phi, \text{All}, \text{work})$  P(no|A||, work, and) P(no|A||, work, and, no)

P(makes|and,no,play)P(Jack|no,play,makes)P(a|play,makes,Jack)P(dull|makes,Jack,a)P(boy|Jack,a,dull)

where, for each factor,

$$[\text{and}|\phi, \text{All}, \text{work}) = \frac{(\text{All work})}{N_{(\text{All work})}}$$

#### Example, a 4-gram model

e: All work and no play makes Jack a dull boy

$$\begin{split} P(e) &= P(\texttt{All}|\phi, \phi, \phi) \; P(\texttt{work}|\phi, \phi, \texttt{All}) \; P(\texttt{and}|\phi, \texttt{All}, \texttt{work}) \\ &\quad P(\texttt{no}|\texttt{All}, \texttt{work}, \texttt{and}) \; P(\texttt{play}|\texttt{work}, \texttt{and}, \texttt{no}) \\ &\quad P(\texttt{makes}|\texttt{and}, \texttt{no}, \texttt{play}) P(\texttt{Jack}|\texttt{no}, \texttt{play}, \texttt{makes}) \\ &\quad P(\texttt{a}|\texttt{play}, \texttt{makes}, \texttt{Jack}) P(\texttt{dull}|\texttt{makes}, \texttt{Jack}, \texttt{a}) \\ &\quad P(\texttt{boy}|\texttt{Jack}, \texttt{a}, \texttt{dull}) \end{split}$$

#### Example, a 4-gram model

e: All work and no play makes Jack a dull boy

$$\begin{split} P(e) &= P(\texttt{All}|\phi, \phi, \phi) \; P(\texttt{work}|\phi, \phi, \texttt{All}) \; P(\texttt{and}|\phi, \texttt{All}, \texttt{work}) \\ &\quad P(\texttt{no}|\texttt{All}, \texttt{work}, \texttt{and}) \; P(\texttt{play}|\texttt{work}, \texttt{and}, \texttt{no}) \\ &\quad P(\texttt{makes}|\texttt{and}, \texttt{no}, \texttt{play}) P(\texttt{Jack}|\texttt{no}, \texttt{play}, \texttt{makes}) \\ &\quad P(\texttt{a}|\texttt{play}, \texttt{makes}, \texttt{Jack}) P(\texttt{dull}|\texttt{makes}, \texttt{Jack}, \texttt{a}) \\ &\quad P(\texttt{boy}|\texttt{Jack}, \texttt{a}, \texttt{dull}) \end{split}$$

e: All work and no play makes Jack a dull boy

$$\begin{split} P(e) &= P(\text{All}|\phi, \phi, \phi) \ P(\text{work}|\phi, \phi, \text{All}) \ P(\text{and}|\phi, \text{All}, \text{work}) \\ P(\text{no}|\text{All}, \text{work}, \text{and}) \ P(\text{play}|\text{work}, \text{and}, \text{no}) \\ P(\text{makes}|\text{and}, \text{no}, \text{play}) P(\text{Jack}|\text{no}, \text{play}, \text{makes}) \\ P(\text{a}|\text{play}, \text{makes}, \text{Jack}) P(\text{dull}|\text{makes}, \text{Jack}, \text{a}) \\ P(\text{boy}|\text{Jack}, \text{a}, \text{dull}) \end{split}$$

#### Example, a 4-gram model

e: All work and no play makes Jack a dull boy

$$\begin{split} P(e) &= P(\texttt{All}|\phi, \phi, \phi) \; P(\texttt{work}|\phi, \phi, \texttt{All}) \; P(\texttt{and}|\phi, \texttt{All}, \texttt{work}) \\ &\quad P(\texttt{no}|\texttt{All}, \texttt{work}, \texttt{and}) \; P(\texttt{play}|\texttt{work}, \texttt{and}, \texttt{no}) \\ &\quad P(\texttt{makes}|\texttt{and}, \texttt{no}, \texttt{play}) P(\texttt{Jack}|\texttt{no}, \texttt{play}, \texttt{makes}) \\ &\quad P(\texttt{a}|\texttt{play}, \texttt{makes}, \texttt{Jack}) P(\texttt{dull}|\texttt{makes}, \texttt{Jack}, \texttt{a}) \\ &\quad P(\texttt{boy}|\texttt{Jack}, \texttt{a}, \texttt{dull}) \end{split}$$

e: All work and no play makes Jack a dull boy

$$\begin{split} P(e) &= P(\text{All}|\phi, \phi, \phi) \ P(\text{work}|\phi, \phi, \text{All}) \ P(\text{and}|\phi, \text{All}, \text{work}) \\ P(\text{no}|\text{All}, \text{work}, \text{and}) \ P(\text{play}|\text{work}, \text{and}, \text{no}) \\ P(\text{makes}|\text{and}, \text{no}, \text{play}) P(\text{Jack}|\text{no}, \text{play}, \text{makes}) \\ P(\text{a}|\text{play}, \text{makes}, \text{Jack}) P(\text{dull}|\text{makes}, \text{Jack}, \text{a}) \\ P(\text{boy}|\text{Jack}, \text{a}, \text{dull}) \end{split}$$

e: All work and no play makes Jack a dull boy

$$\begin{split} P(e) &= P(\text{All}|\phi, \phi, \phi) \; P(\text{work}|\phi, \phi, \text{All}) \; P(\text{and}|\phi, \text{All, work}) \\ &\quad P(\text{no}|\text{All, work, and}) \; P(\text{play}|\text{work, and, no}) \\ &\quad P(\text{makes}|\text{and, no, play}) P(\text{Jack}|\text{no, play, makes}) \\ &\quad P(\text{a}|\text{play, makes, Jack}) P(\text{dull}|\text{makes, Jack, a}) \\ &\quad P(\text{boy}|\text{Jack, a, dull}) \end{split}$$

e: All work and no play makes Jack a dull boy

$$\begin{split} P(e) &= P(\texttt{All}|\phi, \phi, \phi) \; P(\texttt{work}|\phi, \phi, \texttt{All}) \; P(\texttt{and}|\phi, \texttt{All}, \texttt{work}) \\ &= P(\texttt{no}|\texttt{All}, \texttt{work}, \texttt{and}) \; P(\texttt{play}|\texttt{work}, \texttt{and}, \texttt{no}) \\ &= P(\texttt{makes}|\texttt{and}, \texttt{no}, \texttt{play}) P(\texttt{Jack}|\texttt{no}, \texttt{play}, \texttt{makes}) \\ &= P(\texttt{a}|\texttt{play}, \texttt{makes}, \texttt{Jack}) P(\texttt{dull}|\texttt{makes}, \texttt{Jack}, \texttt{a}) \\ &= P(\texttt{boy}|\texttt{Jack}, \texttt{a}, \texttt{dull}) \end{split}$$

e: All work and no play makes Jack a dull boy

$$\begin{split} P(e) &= P(\texttt{All}|\phi, \phi, \phi) \; P(\texttt{work}|\phi, \phi, \texttt{All}) \; \frac{P(\texttt{and}|\phi, \texttt{All}, \texttt{work})}{P(\texttt{no}|\texttt{All}, \texttt{work}, \texttt{and}) \; P(\texttt{play}|\texttt{work}, \texttt{and}, \texttt{no})} \\ &= P(\texttt{makes}|\texttt{and}, \texttt{no}, \texttt{play}) P(\texttt{Jack}|\texttt{no}, \texttt{play}, \texttt{makes}) \\ &= P(\texttt{a}|\texttt{play}, \texttt{makes}, \texttt{Jack}) P(\texttt{dull}|\texttt{makes}, \texttt{Jack}, \texttt{a}) \\ &= P(\texttt{boy}|\texttt{Jack}, \texttt{a}, \texttt{dull}) \end{split}$$

- Long-range dependencies are lost.
- Still, some *n*-grams can be not observed in the corpus.

# Solution

Smoothing techniques:

• Linear interpolation.

$$P(\texttt{and}|\texttt{All},\texttt{work}) = -\frac{N_{(\texttt{All},\texttt{work},\texttt{and})}}{N_{(\texttt{All},\texttt{work})}} + \lambda_2 \frac{N_{(\texttt{work},\texttt{and})}}{N_{(\texttt{work})}} + \lambda_1 \frac{N_{(\texttt{and})}}{N_{words}} + \lambda_0$$

- Long-range dependencies are lost.
- Still, some *n*-grams can be not observed in the corpus.

# Solution

Smoothing techniques:

- Linear interpolation.
- Back-off models.

- Long-range dependencies are lost.
- Still, some *n*-grams can be not observed in the corpus.

# Solution

Smoothing techniques:

• Linear interpolation.

$$P(\texttt{and}|\texttt{All},\texttt{work}) = -\frac{N_{(\texttt{All},\texttt{work},\texttt{and})}}{N_{(\texttt{All},\texttt{work})}} + \lambda_2 \frac{N_{(\texttt{work},\texttt{and})}}{N_{(\texttt{work})}} + \lambda_1 \frac{N_{(\texttt{and})}}{N_{\texttt{words}}} + \lambda_0$$

- Long-range dependencies are lost.
- Still, some *n*-grams can be not observed in the corpus.

# Solution

Smoothing techniques:

• Linear interpolation.

$$P(\texttt{and}|\texttt{All},\texttt{work}) = \lambda_3 \frac{N_{(\texttt{All},\texttt{work},\texttt{and})}}{N_{(\texttt{All},\texttt{work})}} + \lambda_2 \frac{N_{(\texttt{work},\texttt{and})}}{N_{(\texttt{work})}} + \lambda_1 \frac{N_{(\texttt{and})}}{N_{words}} + \lambda_0$$

#### 

cluster:/home/quest/corpus/lm> ls -lkh

-rw-r--r-- 1 emt ia 507M mar 3 15:28 europarl.lm -rw-r--r-- 1 emt ia 50M mar 3 15:29 nc.lm -rw-r--r-- 1 emt ia 3,1G mar 3 15:33 un.lm

cluster:/home/quest/corpus/lm> wc -l

15,181,883 europarl.lm 1,735,721 nc.lm 82,504,380 un.lm

The language model P(e)

cluster:/home/quest/corpus/lm> more nc.lm

\data\ ngram 1=655770 ngram 2=11425501 ngram 3=10824125 ngram 4=13037011 ngram 5=12127575 \1-grams: -3.142546 ! -1.415594-1.978775 " -0.9078496 -4.266428 # -0.2729652-3.806078 \$ -0.3918373 -3.199419 % -1.139753 -3.613416 & -0.6046973 -2.712332 ' -0.6271471 -2.268107 ( -0.6895114

#### The language model P(e)

\2-grams: -1.08232 concierto , -1.093977 concierto . -0.2378127 -1.747908 concierto ad -1.748422 concierto cobraria -0.8927398 concierto de -1.744176 concierto europeo -1.740879 concierto internacional -1.635606 concierto para -1.744787 concierto regional

• • •

\5-grams: -0.8890668 no son los unicos culpables -1.396196 no son los unicos problemas -0.7550655 no son los unicos que -1.240193 no son los unicos responsables

#### Language model: keep in mind

- Statistical LMs estimate the probability of a sentence from its n-gram frequency counts in a monolingual corpus.
- Within an SMT system, it contributes to select fluent sentences in the target language.
- Smoothing techniques are used so that not frequent translations are not discarded beforehand.

The translation model P(f|e)

### Translation model

$$T(f) = \hat{e} = \operatorname{argmax}_{e} P(e) P(f|e)$$

Estimation of the lexical correspondence between languages.

#### How can be P(f|e) characterised?



The translation model P(f|e)

### Translation model

$$T(f) = \hat{e} = \operatorname{argmax}_{e} P(e) P(f|e)$$

Estimation of the lexical correspondence between languages.

#### How can be P(f|e) characterised?



The translation model P(f|e)

### Translation model

$$T(f) = \hat{e} = \operatorname{argmax}_{e} P(e) P(f|e)$$

Estimation of the lexical correspondence between languages.

#### How can be P(f|e) characterised?



The translation model P(f|e)



One should at least model for *each word* in the source language:

- Its translation,
- the number of necessary words in the target language,
- the position of the translation within the sentence,
- and, besides, the number of words that need to be generated from scratch.

### Word-based models: the IBM models

They characterise P(f|e) with 4 parameters: t, n, d and  $p_1$ .

- Lexical probability t t(Quan|When): the prob. that Quan translates into When.
- Fertility n
  n(3|tornes): the prob. that tornes generates 3 words.

### Word-based models: the IBM models

They characterise P(f|e) with 4 parameters: t, n, d and  $p_1$ .

• Distortion d

d(j|i, m, n): the prob. that the word in the *j* position generates a word in the *i* position. *m* and *n* are the length of the source and target sentences.

Probability p1
 p(you|NULL): the prob. that the spurious word you is generated (from NULL).











#### Word-based models: the IBM models

How can be t, n, d and  $p_1$  estimated?

# • Statistical model $\Rightarrow$ counts in a (huge) corpus!

#### But...

• Corpora are aligned at sentence level, not at word level.

#### Solutions

- Pay someone to align 2 milion sentences word by word.
- Estimate word alignments together with the parameters.

#### Word-based models: the IBM models

How can be t, n, d and  $p_1$  estimated?

• Statistical model  $\Rightarrow$  counts in a (huge) corpus!

#### But...

• Corpora are aligned at sentence level, not at word level.

### Solutions

- Pay someone to align 2 milion sentences word by word.
- Estimate word alignments together with the parameters.

#### Word-based models: the IBM models

How can be t, n, d and  $p_1$  estimated?

• Statistical model  $\Rightarrow$  counts in a (huge) corpus!

#### But...

• Corpora are aligned at sentence level, not at word level.

### Solutions

- Pay someone to align 2 milion sentences word by word.
- Estimate word alignments together with the parameters.

The translation model P(f|e)

#### **Expectation-Maximisation algorithm**


The translation model P(f|e)

### **Expectation-Maximisation algorithm**



The translation model P(f|e)

### **Expectation-Maximisation algorithm**



### Alignment's asymmetry

The definitions in IBM models make the alignments asymmetric

• each target word corresponds to only one source word, but the opposite is not true due to the definition of fertility.



### Alignment's asymmetry

The definitions in IBM models make the alignments asymmetric

• each target word corresponds to only one source word, but the opposite is not true due to the definition of fertility.



### Graphically:



Catalan to English

### Graphically:



English to Catalan

Alignment symmetrisation

• Intersection: high-confidence, high precision.



Catalan to English  $\bigcap$  English to Catalan

Alignment symmetrisation

• Union: lower confidence, high recall.



Catalan to English  $\bigcup$  English to Catalan

#### 

cluster:/home/moses/giza.en-es> zmore en-es.A3.final.gz

```
# Sentence pair (1) source length 5 target length 4 alignment score: 0.00015062
resumption of the session
NULL ({ }) reanudacion ({ 1 }) del ({ 2 3 }) periodo ({ }) de ({ }) sesiones ({ 4 })
# Sentence pair (2) source length 33 target length 40 alignment score: 3.3682e-61
```

```
# Sentence pair (2) source length 33 target length 40 alignment score: 3.3682e-61 i declare resumed the session of the european parliament adjourned on friday 17 december 1999, and i would like once again to wish you a happy new year in the hope that you enjoyed a pleasant festive period.
NULL ({ 31 } declare ({ 1 }) reanudado ({ 2 3 }) el ({ 4 }) periodo ({ }) de ({ }) sesiones ({ 5 }) del ({ 6 7 }) parlamento ({ 9 }) europeo ({ 8 }), ({ }) interrumpido ({ 10 }) el ({ }) viernes ({ 12 14 }) 17 ({ 11 13 }) de ({ }) de ({ }) de ({ 16 }) y ({ 16 }) y ({ 17 }) reiretor ({ 21 }) a ({ 23 }) sus ({ 30 }) senorias ({ }) mi ({ 18 }) deseo ({ 24 }) de ({ }) que ({ 33 }) hayan ({ 25 34 35 }) tenido ({ }) unas ({ 19 20 }) buenas ({ 26 36 }) vacaciones ({ 22 27 28 29 32 37 38 39 }). ({ 40 })
```

#### 

cluster:/home/moses/giza.es-en> zmore es-en.A3.final.gz

```
# Sentence pair (1) source length 4 target length 5 alignment score: 1.08865e-07 reanudacion del periodo de sesiones
NULL (\{4\}) resumption (\{1\}) of (\{2\}) the (\{\}) session (\{35\})
```

# Sentence pair (2) source length 40 target length 33 alignment score: 1.88268e-50 declaro reanudado el periodo de sesiones del parlamento europeo , interrumpido el viernes 17 de diciembre pasado , y reitero a sus senorias mi deseo de que hayan tenido unas buenas vacaciones . NULL ({ 5 10 }) i ({ }) declare ({ 1 }) resumed ({ 2 }) the ({ 3 }) session ({ 4 6 }) of ({ 7 }) the ({ }) european ({ 9 }) parliament ({ 8 12 }) adjourned ({ 11 }) on ({ 15 }) friday ({ 13 }) 17 ({ 14 }) december ({ 16 17 }) 1999 ({ }) , ({ 18 }) and ({ 19 }) i ({ }) would ({ }) like ({ }) once ({ }) again ({ }) to ({ 21 }) wish ({ }) you ({ }) a ({ }) happy ({ }) new ({ }) year ({ }) in ({ 26 }) the ({ }) hope ({ }) ) that ({ 27 }) you ({ }) enjoyed ({ 20 }) a ({ }) pleasant ({ 22 23 24 25 28 29 }) festive ({ 30 31 32 }) period ({ }) .

cluster:/home/moses/model> more aligned.grow-diag-final

0-0 1-1 1-2 2-3 4-3

0-0 0-1 1-1 1-2 2-3 3-4 5-4 6-5 6-6 8-7 7-8 11-8 10-9 13-10 14-10 12-11 13-12 12-13 15-14 17-15 18-16 23-17 19-20 20-22 24-23 21-29 26-32 27-33 27-34 30-35 28-36 31-36 29-37 30-37 31-37 31-38 32-39

The translation model P(f|e)

cluster:/home/moses/model> more lex.e2f

```
tuneles tunnels 0.7500000
tuneles transit 0.2000000
estructuralmente weak 1.0000000
estructuralmente structurally 0.5000000
destruido had 0.0454545
para tunnels 0.2500000
sean transit 0.2000000
transito transit 0.6000000
...
```

cluster:/home/moses/model> more lex.f2e

```
tunnels tuneles 0.7500000
transit tuneles 0.2500000
weak estructuralmente 0.5000000
structurally estructuralmente 0.5000000
...
```

The translation model P(f|e)





The translation model P(f|e)

#### From Word-based to Phrase-based models

f: En David llegeix el llibre nou.

#### From Word-based to Phrase-based models

f: En David llegeix el llibre nou. e:  $\phi$ 

#### From Word-based to Phrase-based models

f: En David llegeix el llibre nou.

e: David

From Word-based to Phrase-based models

f: En David llegeix el llibre nou.

e: David reads

- f: En David llegeix el llibre nou.
- e: David reads the

From Word-based to Phrase-based models

f: En David llegeix el llibre nou.

e: David reads the book

From Word-based to Phrase-based models

f: En David llegeix el llibre nou.

e: David reads the book new.

From Word-based to Phrase-based models

f: En David llegeix el llibre nou.

e: David reads the book new.  $\sim$ 

From Word-based to Phrase-based models

f: En David llegeix el llibre nou.

e: David reads the new book.

- f: En David llegeix el llibre nou.
- e: David reads the new book.
- f: En David llegeix el llibre de nou.

#### From Word-based to Phrase-based models

- f: En David llegeix el llibre nou.
- e: David reads the new book.
- f: En David llegeix el llibre de nou.

e: **ø** 

- f: En David llegeix el llibre nou.
- e: David reads the new book.
- f: En David llegeix el llibre de nou.
- e: David

- f: En David llegeix el llibre nou.
- e: David reads the new book.
- f: En David llegeix el llibre de nou.
- e: David reads

- f: En David llegeix el llibre nou.
- e: David reads the new book. 🗸
- f: En David llegeix el llibre de nou.
- e: David reads the

- f: En David llegeix el llibre nou.
- e: David reads the new book. 🗸
- f: En David llegeix el llibre de nou.
- e: David reads the book

- f: En David llegeix el llibre nou.
- e: David reads the new book. 🗸
- f: En David llegeix el llibre de nou.
- e: David reads the book of

- f: En David llegeix el llibre nou.
- e: David reads the new book. 🗸
- f: En David llegeix el llibre de nou.
- e: David reads the book of new.

#### From Word-based to Phrase-based models

f: En David llegeix el llibre nou.

e: David reads the new book. 🗸

f: En David llegeix el llibre de nou.

e: David reads the book of new. 🗡

#### From Word-based to Phrase-based models

f: En David llegeix el llibre nou.

e: David reads the new book. 🗸

f: En David llegeix el llibre de nou. e: David reads the book of new. Xe:  $\phi$ 

- f: En David llegeix el llibre nou.
- e: David reads the new book. 🗸
- f: En David llegeix el llibre de nou.
- e: David reads the book of new. 🗡 e: David

- f: En David llegeix el llibre nou.
- e: David reads the new book. 🗸
- f: En David llegeix el llibre de nou.
- e: David reads the book of new. 🗡 e: David reads

#### From Word-based to Phrase-based models

f: En David llegeix el llibre nou.

e: David reads the new book. 🗸

f: En David llegeix el llibre de nou.

e: David reads the book of new. X e: David reads the

#### From Word-based to Phrase-based models

f: En David llegeix el llibre nou.

e: David reads the new book. 🗸

f: En David llegeix el llibre de nou.

e: David reads the book of new. X e: David reads the book
### From Word-based to Phrase-based models

f: En David llegeix el llibre nou.

e: David reads the new book. 🗸

f: En David llegeix el llibre de nou.

e: David reads the book of new. 🗡 e: David reads the book <mark>again</mark>.

### From Word-based to Phrase-based models

f: En David llegeix el llibre nou.

e: David reads the new book.

f: En David llegeix el llibre de nou.

e: David reads the book of new. imes e: David reads the book again.  $\checkmark$ 

### From Word-based to Phrase-based models

f: En David llegeix el llibre nou.

e: David reads the new book. 🗸

f: En David llegeix el llibre de nou.

e: David reads the book of new. 🗡 e: David reads the book again. 🗸

- Some sequences of words usually translate together.
- Approach: take sequences (phrases) as translation units.

## What can be achieved with phrase-based models (as compared to word-based models)

- Allow to translate from several to several words and not only from one to several.
- Some local and short range context is used.
- Idioms can be catched.

The translation model P(f|e)



With the new translation units, P(f|e) can be obtained following the same strategy as for word-based models with few modifications:

- Segment source sentence in phrases.
- 2 Translate each phrase into the target language.
- 3 Reorder the output.

The translation model P(f|e)



With the new translation units, P(f|e) can be obtained following the same strategy as for word-based models with few modifications:

- Segment source sentence in phrases.
- 2 Translate each phrase into the target language.
- Reorder the output.

The translation model P(f|e)



With the new translation units, P(f|e) can be obtained following the same strategy as for word-based models with few modifications:

- Segment source sentence in phrases.
- **2** Translate each phrase into the target language.
- 8 Reorder the output.

The translation model P(f|e)



### But...

• Alignments need to be done at phrase level

## Options

- Calculate phrase-to-phrase alignments  $\Rightarrow$  hard!
- Obtain phrase alignments from word alignments  $\Rightarrow$  how?

Questions to answer:

- How do we obtain phrase alignments from word alignments?
- And, by the way, what's exactly a phrase?!

A **phrase is** a sequence of words consistent with word alignment. That is, no word is aligned to a word outside the phrase. But a phrase **is not** necessarily a linguistic element.

Questions to answer:

- How do we obtain phrase alignments from word alignments?
- And, by the way, what's exactly a phrase?!

A **phrase is** a sequence of words consistent with word alignment. That is, no word is aligned to a word outside the phrase. But a phrase **is not** necessarily a linguistic element.

Questions to answer:

- How do we obtain phrase alignments from word alignments?
- And, by the way, what's exactly a phrase?!

A **phrase is** a sequence of words consistent with word alignment. That is, no word is aligned to a word outside the phrase. But a phrase **is not** necessarily a linguistic element.

Questions to answer:

- How do we obtain phrase alignments from word alignments?
- And, by the way, what's exactly a phrase?!

A **phrase is** a sequence of words consistent with word alignment. That is, no word is aligned to a word outside the phrase. But a phrase **is not** necessarily a linguistic element.<sup>1</sup>

### Phrase extraction through an example:



(Quan tornes, When are you coming back)

### Phrase extraction through an example:



(Quan tornes, When are you coming back)

### Phrase extraction through an example:



### The translation model P(f|e)

### 

### The translation model P(f|e)

## Quan tornesacasa?When<br/>are<br/>youIIIintersectionIIIorningIIIback<br/>homeIII?III

### The translation model P(f|e)

## IntersectionQuan tornes a casa ?When<br/>are<br/>youQuan tornes a casa ?When<br/>are<br/>youImage: State of the state of

### The translation model P(f|e)

# Quan tornesacasa?When<br/>are<br/>you<br/>comingIIIare<br/>you<br/>comingIIIback<br/>homeIII?III

### The translation model P(f|e)

# Quan tornesacasa?When<br/>are<br/>you<br/>comingIIIare<br/>you<br/>comingIIIback<br/>homeIII?III

### The translation model P(f|e)

# Quan tornesacasa?When<br/>are<br/>youImage: Coming<br/>back<br/>homeImage: Coming<br/>point of the sector of the sector

### The translation model P(f|e)

# Quan tornesacasa?When<br/>are<br/>youImage: Coming<br/>back<br/>homeImage: Coming<br/>point of the comingImage: Coming<br/>point of the coming<br/>point of th

### The translation model P(f|e)

# Quan tornesacasa?When<br/>are<br/>youImage: Coming<br/>back<br/>homeImage: Coming<br/>point of the comingImage: Coming<br/>point of the coming<br/>point of th

### The translation model P(f|e)

# Quan tornesacasa?When<br/>are<br/>you<br/>comingImage: Common tornesImage: Common tornesImage: Common tornesImage: Common tornesback<br/>home<br/>?Image: Common tornesImage: Common tornesImage: Common tornesImage: Common tornescoming<br/>?Image: Common tornesImage: Common tornesImage: Common tornesImage: Common tornes

### The translation model P(f|e)

# Quan tornesacasa?When<br/>are<br/>you<br/>comingImage: Common com

### The translation model P(f|e)

# UnionQuan tornes a casa ?When<br/>are<br/>you<br/>coming<br/>back<br/>nomeImage: Casa (Casa) (Cas

(Quan, When) (Quan tornes, When are) (Quan tornes, When are you coming) (Quan tornes, When are you coming back) (Quan tornes a casa, When are you coming back home) ... (tornes a casa ?, are you coming back home ?) (casa, home) (casa ?, home ?) (?, ?) 21 phrases

### The translation model P(f|e)

# Quan tornes a casa ? When are you Image: Second second

(Quan, When) (Quan tornes, When are) (Quan tornes, When are you coming) (Quan tornes, When are you coming back) (Quan tornes a casa, When are you coming back home) .... (tornes a casa ?, are you coming back home ?) (casa, home) (casa ?, home ?) (?, ?) 21 phrases

### The translation model P(f|e)

# UnionQuan tornes a casa ?When<br/>are<br/>you<br/>comingImage: Casa (Casa)When<br/>are<br/>you<br/>comingImage: Casa (Casa)back<br/>home<br/>?Image: Casa (Casa)?Image: Casa (Casa)

(Quan, When) (Quan tornes, When are) (Quan tornes, When are you coming) (Quan tornes, When are you coming back) (Quan tornes a casa, When are you coming back home) .... (tornes a casa ?, are you coming back home ?) (casa, home) (casa ?, home ?) (?, ?) 21 phrases

### The translation model P(f|e)

# Quan tornes a casa ? When are you Image: Second second

(Quan, When) (Quan tornes, When are) (Quan tornes, When are you coming) (Quan tornes, When are you coming back) (Quan tornes a casa, When are you coming back home) .... (tornes a casa ?, are you coming back home ?) (casa, home) (casa ?, home ?) (?, ?) 21 phrases

### The translation model P(f|e)



(Quan, When) (Quan tornes, When are) (Quan tornes, When are you coming) (Quan tornes, When are you coming back) (Quan tornes a casa, When are you coming back home) ... (tornes a casa ?, are you coming back home ?) (casa, home) (casa ?, home ?) (?, ?) 21 phrases

### Phrase extraction

- The number of extracted phrases depends on the symmetrisation method.
  - Intersection: few precise phrases.
  - Union: lots of (less?) precise phrases.
- Usually, neither intersection nor union are used, but something in between.
  - Start from the intersection and add points belonging to the union according to heuristics.

## Phrase extraction

- For each phrase-pair  $(f_i, e_i)$ ,  $P(f_i|e_i)$  is estimated by frequency counts in the parallel corpus.
- The set of possible phrase-pairs conforms the set of translation options.
- The set of phrase-pairs together with their probabilities conform the translation table.

### 

### cluster:/home/moses/model> zmore extract.gz

reanudacion ||| resumption ||| 0-0 reanudacion del ||| resumption of the ||| 0-0 1-1 1-2 reanudacion del periodo de sesiones ||| resumption of the session ||| 0-0 1-1 1-2 2-3 4-3

### cluster:/home/moses/model> zmore extract.inv.gz

```
resumption ||| reanudacion ||| 0-0
resumption of the ||| reanudacion del ||| 0-0 1-1 2-1
resumption of the session ||| reanudacion del periodo de sesiones ||| 0-0 1-1 2-1 3-2 3-4
```

```
cluster:/home/moses/model> zmore extract.o.gz
```

reanudacion ||| resumption ||| mono mono reanudacion del ||| resumption of the ||| mono mono reanudacion del periodo de sesiones ||| resumption of the session ||| mono mono

The translation model P(f|e)

### cluster:/home/moses/model> zmore phrase-table.gz

```
be consistent ||| coherentes ||| 0.0384615 0.146893 0.083333 0.0116792 2.718 ||| 1-0 ||| 26 12

be consistent ||| sean coherentes ||| 0.2 0.00022714 0.083333 0.0116792 2.718 ||| 0-0 1-1 ||| 5 12

be consistent ||| sean consistentes ||| 0.5 0.00014834 0.083333 0.0785852 2.718 ||| 0-0 1-1 ||| 2 12

be consistent ||| ser coherente ||| 0.5 0.00014834 0.0833333 0.0785852 2.718 ||| 0-0 1-1 ||| 4 12

be consistent ||| ser consecuente ||| 1 0.000340072 0.0833333 0.758942 2.718 ||| 0-0 1-1 ||| 4 12

be consistent ||| ser consecuente ||| 1 0.00851083 0.5 0.633285 2.718 ||| 0-0 1-1 ||| 6 12

consistent ||| ser consistente ||| 1 0.0085183 0.5 0.633285 2.718 ||| 0-0 1-1 ||| 6 12

consistent ||| coherente cuando se ||| 1 0.00783857 1 0.329794 2.718 ||| 0-0 1-1 1-2 ||| 1 1

consistent ||| adecuado ||| 0.00512821 0.0112994 0.00671141 0.00909 2.718 ||| 0-0 ||| 195 149

consistent ||| constante ||| 0.137931 0.0282486 0.026456 0.0847458 2.718 ||| 0-0 ||| 29 149

consistent ||| constantes ||| 0.0625 0.0056497 0.00671141 0.047619 2.718 ||| 0-0 ||| 16 149
```

### Translation model: keep in mind

- Statistical TMs estimate the probability of a translation from a parallel aligned corpus.
- Its quality depends on the quality of the obtained word (phrase) alignments.
- Within an SMT system, it contributes to select semantically adequate sentences in the target language.

## SMT, components Decoder

## Decoder

$$T(f) = \hat{e} = \operatorname{argmax}_{e} P(e) P(f|e)$$

Responsible for the search in the space of possible translations.

Given a model (LM+TM+...), the decoder constructs the possible translations and looks for the most probable one.

### In our context, one can find:

- Greedy decoders. Initial hypothesis (word by word translation) refined iteratively using hill-climbing heuristics.
- Beam search decoders.
### SMT, components Decoder

#### Decoder

$$T(f) = \hat{e} = \operatorname{argmax}_{e} P(e) P(f|e)$$

Responsible for the search in the space of possible translations.

Given a model (LM+TM+...), the decoder constructs the possible translations and looks for the most probable one.

In our context, one can find:

- Greedy decoders. Initial hypothesis (word by word translation) refined iteratively using hill-climbing heuristics.
- Beam search decoders.

### SMT, components Decoder

#### Decoder

$$T(f) = \hat{e} = \operatorname{argmax}_{e} P(e) P(f|e)$$

Responsible for the search in the space of possible translations.

Given a model (LM+TM+...), the decoder constructs the possible translations and looks for the most probable one.

In our context, one can find:

- Greedy decoders. Initial hypothesis (word by word translation) refined iteratively using hill-climbing heuristics.
- Beam search decoders. Let's see..

A beam-search decoder

#### Core algorithm



A beam-search decoder

Example: Quan tornes a casa

• Translation options:

```
(Quan, When)
(Quan tornes, When are you coming back)
(Quan tornes a casa, When are you coming back home)
(tornes, come back)
(tornes a casa, come back home)
(a casa, home)
```

A beam-search decoder

Example: Quan tornes a casa

• Translation options:

```
(Quan, When)
(Quan tornes, When are you coming back)
(Quan tornes a casa, When are you coming back home)
(tornes, come back)
(tornes a casa, come back home)
(a casa, home)
```

• Notation for hypotheses in construction:

Constructed sentence so far:come backSource words already translated:- x - -

A beam-search decoder

Example: Quan tornes a casa

• Translation options:

```
(Quan, When)
(Quan tornes, When are you coming back)
(Quan tornes a casa, When are you coming back home)
(tornes, come back)
(tornes a casa, come back home)
(a casa, home)
```

• Notation for hypotheses in construction:

Constructed sentence so far: come back Source words already translated: - x - -

A beam-search decoder

Example: Quan tornes a casa

• Translation options:

```
(Quan, When)
(Quan tornes, When are you coming back)
(Quan tornes a casa, When are you coming back home)
(tornes, come back)
(tornes a casa, come back home)
(a casa, home)
```

Initial hypothesis

Constructed sentence so far: Source words already translated:

A beam-search decoder

φ











A beam-search decoder

#### Exhaustive search

• As a result, one should have an estimation of the cost of each hypothesis, being the lowest cost one the best translation.

But..

• The number of hypotheses is exponential with the number of source words.

(30 words sentence  $\Rightarrow 2^{30} = 1,073,741,824$  hypotheses!)

#### Solution

- Optimise the search by:
  - Hypotheses recombination
  - Beam search and pruning

A beam-search decoder

#### Exhaustive search

• As a result, one should have an estimation of the cost of each hypothesis, being the lowest cost one the best translation.

But...

• The number of hypotheses is exponential with the number of source words.

(30 words sentence  $\Rightarrow 2^{30} = 1,073,741,824$  hypotheses!)

Solution

- Optimise the search by:
  - Hypotheses recombination
  - Beam search and pruning

A beam-search decoder

#### Exhaustive search

• As a result, one should have an estimation of the cost of each hypothesis, being the lowest cost one the best translation.

But...

• The number of hypotheses is exponential with the number of source words.

(30 words sentence  $\Rightarrow 2^{30} = 1,073,741,824$  hypotheses!)

#### Solution

- Optimise the search by:
  - Hypotheses recombination
  - Beam search and pruning



#### Hypotheses recombination

Combine hypotheses with the same source words translated, keep that with a lower cost.

- Risk-free operation. The lowest cost translation is still there.
- But the space of hypothesis is not reduced enough.

### Hypotheses recombination

Combine hypotheses with the same source words translated, keep that with a lower cost.

- Risk-free operation. The lowest cost translation is still there.
- But the space of hypothesis is not reduced enough.

### Hypotheses recombination

Combine hypotheses with the same source words translated, keep that with a lower cost.

$$\begin{array}{ccc} \text{When} \mid \texttt{come\_back\_home} & \longleftrightarrow & \text{When} \mid \texttt{come\_back} \mid \texttt{home} \\ & \times \times \times & & \times \times \times \end{array}$$

- Risk-free operation. The lowest cost translation is still there.
- But the space of hypothesis is not reduced enough.

#### Beam search and pruning (at last!)

Compare hypotheses with the same number of translated source words and prune out the inferior ones.

What is an inferior hypothesis?

- The quality of a hypothesis is given by the cost so far and by an estimation of the future cost.
- Future cost estimations are only approximate, so the pruning is not risk-free.

### Beam search and pruning (at last!)

Strategy:

- Define a beam size (by threshold or number of hypotheses).
- Distribute the hypotheses being generated in stacks according to the number of translated source words, for instance.
- Prune out the hypotheses falling outside the beam.
- The hypotheses to be pruned are those with a higher (current + future) cost.

Decoding: keep in mind

- Standard SMT decoders translate the sentences from left to right by expanding hypotheses.
- Beam search decoding is one of the most efficient approach.
- But, the search is only approximate, so, the best translation can be lost if one restricts the search space too much.

## 1 Introduction

## 2 Basics

3 Components

4 The log-linear model



Maximum likelihood (ML)

$$\hat{e} = \operatorname{argmax}_{e} P(e|f) = \operatorname{argmax}_{e} P(e) P(f|e)$$

Maximum entropy (ME)

$$\hat{e} = \operatorname{argmax}_{e} P(e|f) = \operatorname{argmax}_{e} \exp\left\{\sum \lambda_{m} h_{m}(f, e)\right\}$$

$$\hat{e} = \operatorname{argmax}_{e} \log P(e|f) = \operatorname{argmax}_{e} \sum \lambda_{m} h_{m}(f, e)$$

Log-linear model

Maximum likelihood (ML)

$$\hat{e} = \operatorname{argmax}_{e} P(e|f) = \operatorname{argmax}_{e} P(e) P(f|e)$$

Maximum entropy (ME)

$$\hat{e} = \operatorname{argmax}_{e} P(e|f) = \operatorname{argmax}_{e} \exp\left\{\sum \lambda_{m} h_{m}(f, e)\right\}$$

$$\hat{e} = \operatorname{argmax}_{e} \log P(e|f) = \operatorname{argmax}_{e} \sum \lambda_{m} h_{m}(f, e)$$

Maximum likelihood (ML)

$$\hat{e} = \operatorname{argmax}_{e} P(e|f) = \operatorname{argmax}_{e} P(e) P(f|e)$$

Maximum entropy (ME)

$$\hat{e} = \operatorname{argmax}_{e} P(e|f) = \operatorname{argmax}_{e} \exp\left\{\sum \lambda_{m} h_{m}(f, e)\right\}$$

$$\hat{e} = \operatorname{argmax}_{e} \log P(e|f) = \operatorname{argmax}_{e} \sum \lambda_{m} h_{m}(f, e)$$
  
Log-linear model

Maximum likelihood (ML)

$$\hat{e} = \operatorname{argmax}_{e} P(e|f) = \operatorname{argmax}_{e} P(e) P(f|e)$$

Maximum entropy (ME)

 $\hat{e} = \operatorname{argmax}_{e} \log P(e|f) = \operatorname{argmax}_{e} \sum \lambda_{m} h_{m}(f, e)$ Log-linear model with  $h_{1}(f, e) = \log P(e), \ h_{2}(f, e) = \log P(f|e), \ \text{and} \ \lambda_{1} = \lambda_{2} = 1$  $\Rightarrow \text{Maximum likelihood model}$ 

## What can achieved with the log-linear model (as compared to maximum likelihood model)

- Extra features  $h_m$  can be easily added...
- ... but their weight  $\lambda_m$  must be somehow determined.
- Different knowledge sources can be used.

#### State of the art feature functions

Eight features are usually used: P(e), P(f|e), P(e|f), lex(f|e), lex(e|f), ph(e), w(e) and  $P_d(e, f)$ .

- Language model P(e)
   P(e): Language model probability as in ML model.
- Translation model P(f|e)
   P(f|e): Translation model probability as in ML model.
- Translation model P(e|f)
   P(e|f): Inverse translation model probability to be added to the generative one.

#### State of the art feature functions

Eight features are usually used: P(e), P(f|e), P(e|f), lex(f|e), lex(e|f), ph(e), w(e) and  $P_d(e, f)$ .

- Translation model lex(f|e)lex(f|e): Lexical translation model probability.
- Translation model lex(e|f)
   lex(e|f): Inverse lexical translation model probability.
- Phrase penalty ph(e)
   ph(e): A constant cost per produced phrase.

#### State of the art feature functions

Eight features are usually used: P(e), P(f|e), P(e|f), lex(f|e), lex(e|f), ph(e), w(e) and  $P_d(e, f)$ .

- Word penalty w(e)
   w(e): A constant cost per produced word.
- Distortion P<sub>d</sub>(e, f)
   P<sub>d</sub>(ini<sub>phrase<sub>i</sub></sub>, end<sub>phrase<sub>i-1</sub></sub>): Relative distortion probability distribution. A simple distortion model:
   P<sub>d</sub>(ini<sub>phrase<sub>i</sub></sub>, end<sub>phrase<sub>i-1</sub></sub>) = α|ini<sub>phrase<sub>i</sub></sub> end<sub>phrase<sub>i-1</sub></sub> 1|

#### SMT, components The translation model P(f|e)

#### 

cluster:/home/moses/model> zmore phrase-table.gz

```
be consistent ||| coherentes ||| 0.0384615 0.146893 0.083333 0.0116792 2.718 ||| 1-0 ||| 26 12

be consistent ||| sean coherentes ||| 0.2 0.00022714 0.0833333 0.0786835 2.718 ||| 0-0 1-1 ||| 5 12

be consistent ||| sean consistentes ||| 0.5 0.00014834 0.0833333 0.0786835 2.718 ||| 0-0 1-1 ||| 2 12

be consistent ||| ser coherente ||| 0.5 0.00014834 0.0833333 0.0786835 2.718 ||| 0-0 1-1 ||| 4 12

be consistent ||| ser consecuente ||| 1 0.00034072 0.083333 0.758942 2.718 ||| 0-0 1-1 ||| 4 12

be consistent ||| ser consistente ||| 1 0.00034077 0.083333 0.758942 2.718 ||| 0-0 1-1 ||| 6 12

consistent ||| ser consistente ||| 1 0.000783857 1 0.329794 2.718 ||| 0-0 1-1 ||| 6 12

consistent ||| adecuado ||| 0.00512821 0.0112994 0.00671141 0.009009 2.718 ||| 0-0 ||| 195 149

consistent ||| coherenta ||| 0.137931 0.0282486 0.0268456 0.0847458 2.718 ||| 0-0 ||| 29 149

consistent ||| constante ||| 0.033333 0.0112994 0.00671141 0.007619 2.718 ||| 0-0 ||| 60 149

consistent ||| constantes ||| 0.0625 0.0056497 0.00671141 0.047619 2.718 ||| 0-0 ||| 61 149
```

. . .

Digression: lexicalised reordering or distortion

#### State of the art?

Software such as Moses makes easy the incorporation of more sophisticated reordering.

From a **distance-based** reordering (1 feature)

to include orientation information in a **lexicalised** reordering. (3-6 features)

## SMT, the log-linear model

Digression: lexicalised reordering or distortion

#### From where and how can one learn reorders?



(are, tornes, monotone)

## SMT, the log-linear model

Digression: lexicalised reordering or distortion

#### From where and how can one learn reorders?



(coming back, tornes, swap)

## SMT, the log-linear model

Digression: lexicalised reordering or distortion

From where and how can one learn reorders?



(home ?, casa ?, discontinuous)
Digression: lexicalised reordering or distortion

3 new features estimated by frequency counts:  $P_{\rm monotone}$ ,  $P_{\rm swap}$  and  $P_{\rm discontinuous}$  (6 when bidirectional).

$$P_{or.}( ext{orientation}|f,e) = rac{count( ext{orientation},e,f)}{\sum_{or.} count( ext{orientation},e,f)}$$

- $\bullet\,$  Sparse statistics of the orientation types  $\rightarrow$  smoothing.
- Several variations.

#### SMT, components The translation model P(f|e)

#### 🖅 In practice,

cluster:/home/moses/model> zmore extract.o.gz

resumption ||| reanudacion ||| mono mono resumption of the ||| reanudacion del ||| mono mono resumption of the session ||| reanudacion del periodo de sesiones ||| mono mono de la union ||| union ' s ||| swap swap competencia de la union ||| union ' s competition ||| swap other ...

cluster:/home/moses/model> zmore reordering-table.wbe-msd-bidirectional-fe.gz

a resumption of the s ||| se reanudara el periodo de s |||  $0.200 \ 0.200 \ 0.600 \ 0.600 \ 0.200 \ 0.200$ resumption of the s ||| reanudacion del periodo de s |||  $0.995 \ 0.002 \ 0.002 \ 0.995 \ 0.002 \ 0.002$ the resumption of the s ||| la continuacion del periodo de s |||  $0.142 \ 0.142 \ 0.714 \ 0.714 \ 0.142 \ 0.142$ the resumption of the s ||| la reanudacion del periodo de s |||  $0.818 \ 0.090 \ 0.090 \ 0.818 \ 0.090 \ 0.090$ 

#### SMT, components The translation model P(f|e)

cluster:/home/moses/model> wc -l \*

```
493,896,818 phrase-table
493,896,818 reordering-table.wbe-msd-bidirectional-fe
```

cluster:/home/moses/model> ls -lkh \*

```
-rw-r--r-- 1 emt ia 57G mar 3 14:01 phrase-table
-rw-r--r-- 1 emt ia 55G mar 3 14:08 reordering-table.wbe-msd-bidirectional-fe
```

#### State of the art feature functions

13 features may be used:

- *P*(*e*);
- P(f|e), P(e|f), lex(f|e), lex(e|f);
- *ph*(*e*), *w*(*e*);
- $P_{mon}(o|e, f)$ ,  $P_{swap}(o|e, f)$ ,  $P_{dis}(o|e, f)$ ,
- $P_{mon}(o|f,e)$ ,  $P_{swap}(o|f,e)$ ,  $P_{dis}(o|f,e)$ .

#### Development training, weights optimisation

• Supervised training: a (small) aligned parallel corpus is used to determine the optimal weights.

$$\hat{e} = \operatorname{argmax}_{e} \log P(e|f) = \operatorname{argmax}_{e} \sum \lambda_{m} h_{m}(f, e)$$

#### Development training, weights optimisation

Strategies

- Generative training. Optimises ME objective function which has a unique optimum. Maximises the likelihood.
- Discriminative training only for feature weights (not models), or purely discriminative for the model as a whole. This way translation performance can be optimised.
- Minimum Error-Rate Training (MERT).

#### Development training, weights optimisation

Strategies

- Generative training. Optimises ME objective function which has a unique optimum. Maximises the likelihood.
- Discriminative training only for feature weights (not models), or purely discriminative for the model as a whole. This way translation performance can be optimised.
- Minimum Error-Rate Training (MERT).

#### Minimum Error-Rate Training

• Approach: Minimise an error function.

But... what's the error of a translation?

- There exist several error measures or metrics.
- Metrics not always correlate with human judgements.
- The quality of the final translation on the metric choosen for the optimisation is shown to improve.
- For the moment, let's say we use BLEU.

(More on MT Evaluation section)

# SMT, the log-linear model Minimum Error-Rate Training (MERT)

#### Minimum Error-Rate Training rough algorithm



The log-linear model

#### Log-linear model: keep in mind

- The log-linear model allows to include several weighted features. State of the art systems use 8 real features.
- The corresponding weights are optimised on a development set, a small aligned parallel corpus.
- An optimisation algorithm such as MERT is appropriate for at most a dozen of features. For more features, purely discriminative learnings should be used.
- For MERT, the choice of the metric that quantifies the error in the translation is an issue.

# 1 Introduction



3 Components

4 The log-linear model

# Beyond standard SMT

- Factored translation models
- Syntactic translation models
- Ongoing research

#### Considering linguistic information in phrase-based models

• Phrase-based log-linear models do not consider linguistic information other than words. This is information should be included.

# Options

- Use syntactic information as pre- or post-process (for reordering or reranking for example).
- Include linguistic information in the model itself.
  - Factored translation models.
  - Syntactic-based translation models.

Factored translation models

#### **Factored translation models**

Extension to phrase-based models where every word is substituted by a vector of factors.

 $(word) \Longrightarrow (word, lemma, PoS, morphology, ...)$ 

The translation is now a combination of pure translation (T) and generation (G) steps:

Factored translation models

#### **Factored translation models**

Extension to phrase-based models where every word is substituted by a vector of factors.

$$(word) \Longrightarrow (word, lemma, PoS, morphology, ...)$$

The translation is now a combination of pure translation (T) and generation (G) steps:

Factored translation models

#### **Factored translation models**

Extension to phrase-based models where every word is substituted by a vector of factors.

$$(word) \Longrightarrow (word, lemma, PoS, morphology, ...)$$

The translation is now a combination of pure translation (T) and generation (G) steps:

$$\begin{array}{ccc} casa_{f} & NN_{f} & fem., plural_{f} & cases_{f} \\ \downarrow \top & \downarrow \top & \downarrow \top \\ house_{e} & NN_{e} & plural_{e} & \xrightarrow{G} \\ houses_{e} \end{array}$$

# What differs in factored translation models (as compared to standard phrase-based models)

- The parallel corpus must be annotated beforehand.
- Extra language models for every factor can also be used.
- Translation steps are accomplished in a similar way.
- Generation steps imply a training only on the target side of the corpus.
- Models corresponding to the different factors and components are combined in a log-linear fashion.

# SMT, beyond standard SMT

Syntactic translation models

#### Syntactic translation models

Incorporate syntax to the source and/or target languages.

#### Approaches

- Syntactic phrase-based based on tree trasducers:
  - Tree-to-string. Build mappings from target parse trees to source strings.
  - String-to-tree. Build mappings from target strings to source parse trees.
  - ► Tree-to-tree. Mappings from parse trees to parse trees.

Syntactic translation models

# Syntactic translation models

Incorporate syntax to the source and/or target languages.

# Approaches

- Synchronous grammar formalism which learns a grammar that can simultaneously generate both trees.
  - ► Syntax-based. Respect linguistic units in translation.
  - Hierarchical phrase-based. Respect phrases in translation.





















Syntactic models ease reordering. An intuitive example:







David reads a new book

#### Hot research topics

Current research on SMT addresses known and new problems.

Some components of the standard phrase-based model are still under study:

- Automatic alignments.
- Language models and smoothing techniques.
- Parameter optimisation.

Complements to a standard system can be added:

- Reordering as a pre-process or post-process.
- Reranking of n-best lists.
- OOV treatment.
- Domain adaptation.

Development of full systems from scratch or modifications to the standard:

- Using machine learning.
- Including linguistic information.
- Hybridation of MT paradigms.
- Or a different strategy:
  - Systems combination.

# SMT, beyond standard SMT Including linguistic information

#### Beyond standard SMT: keep in mind

- Factored models include linguistic information in phrasebased models and are suitable for morphologically rich languages.
- Syntactic models consider somehow syntaxis and are adequate for language pairs with a different structure of the sentences.
- Current research addresses both new models and modifications to the existing ones.

# Part II

# SMT experiments



# 6 Translation system

- Demos
- Software
- Steps

# SMT system Demo: http://demo.statmt.org/

O Moses Online MT Demo - Mo	zilla Firefox					_	
Eitxer Edita Visualitza Historial Adreces d'interès Eines Ajuda							
🖕 🗼 🗸 🕑 🔝 🏫 间 http://demo.statmt.org/index.php				☆ 🗸 😽 Go	ogle		0,
⊘ Disable ✓ Scookies ✓ CCSS ✓ Forms ✓ Images ✓ Information ✓ ③ Miscellaneous	;∽ ∥Outline∽	§ ∎Resize ∽	🥜 Tools 🗸	⊇View Source ✓	⊘Options ∨	×	0 0
S Moses Online MT Demo							~
Moses Machine Translation Demo Source: Try any example to translate. English->German © Show Debug Output Show Alignment Translate Looking to translate a web page? Then click here						The second se	
Translation: Versuchen Sie ein Beispiel zu übersetzen. Help to Improve statistical machine translation! Versuchen Sie ein Beispiel zu übersetzen.	]						~
Fet						4	*

# SMT system Demo: http://cog.hut.fi/smtdemo

Translator demo - Mozilla Firefox	_ D X
<u>E</u> itxer <u>E</u> dita <u>V</u> isualitza Hi <u>s</u> torial Ad <u>r</u> eces d'interès Ei <u>n</u> es Ajuda	
🖕 🧼 🗸 🧭 🔞 http://cog.hut.fi/smtdemo	🛃 🗸 🔍
© Disable ∨ 💩 Cookies ∨ 🔤 CSS ∨ 🔄 Forms ∨ 🔳 Images × 🚯 Information × ⊗Miscellaneous × 🖉 Outline × 👯 Resize × 🤌 Tools × 😥 View	Source 🗸 🖉 Options 🗸 🛷 🕲 🕲
🗑 Translator demo 📫	~
Statistical Machine Translation Demo This page domentrates the idea presented in the pager '5. Virpigi, j. j. Värpige, M. Creatz, M. Sadenieni, Merphology-Aware Statistical Machine Translations Based on Merphs Indexed in an Unsuper- Copendance, Downmark, pp. 012-007, 2007. Try the translator below, or view recent translations. Try an example to translate.	teed Manner. In Proceedings of MT Summit XI,
en->sv (word, Europarl v3) 0 Translate	
Show phrases horizontal C	
This page is maintained by the <u>Commutational Coupline Systems Group</u> at the Aalto University.	
Fet	*



#### Build your own SMT system

- Language model with SRILM. http://wwwspeech.sri.com/projects/srilm/download.html
- Word alignments with GIZA++. http://code.google.com/p/giza-pp/downloads/list
- And everything else with the Moses package. https://github.com/moses-smt/mosesdecoder

# 1. Download and prepare your data

Parallel corpora and some tools can be downloaded for instance from the WMT 2013 web page: http://www.statmt.org/wmt13/translation-task.html

How to construct a baseline system is also explained there:  $\label{eq:http://www.statmt.org/wmt10/baseline.html}$ 

We continue with the Europarl corpus Spanish-to-English.

#### 1. Download and prepare your data (cont'd)

 Tokenise the corpus with WMT10 scripts. (training corpus and development set for MERT)

```
wmt10scripts/tokenizer.perl -l es < eurov4.es-en.NOTOK.es >
eurov4.es-en.TOK.es
wmt10scripts/tokenizer.perl -l en < eurov4.es-en.NOTOK.en >
eurov4.es-en.TOK.en
```

```
wmt10scripts/tokenizer.perl -l es < eurov4.es-en.NOTOK.dev.es >
eurov4.es-en.TOK.dev.es
wmt10scripts/tokenizer.perl -l en < eurov4.es-en.NOTOK.dev.en >
eurov4.es-en.TOK.dev.en
```



- 1. Download and prepare your data (cont'd)
  - Filter out long sentences with Moses scripts. (Important for GIZA++)

bin/moses-scripts/training/clean-corpus-n.perl eurov4.es-en.TOK es
en eurov4.es-en.TOK.clean 1 100

 Lowercase training and development with WMT10 scripts. (Optional but recommended)

```
wmt10scripts/lowercase.perl < eurov4.es-en.TOK.clean.es >
eurov4.es-en.es
wmt10scripts/lowercase.perl < eurov4.es-en.TOK.clean.en >
eurov4.es-en.en
```

#### 2. Build the language model

 Run SRILM on the English part of the parallel corpus or on a monolingual larger one. (tokenise and lowercase in case it is not)

ngram-count -order 5 -interpolate -kndiscount -text eurov4.es-en.en -lm eurov4.en.lm
### SMT system Steps

### 3. Train the translation model

• Use the Moses script train-model.perl This script performs the whole training:

```
train-model.perl -help
```

Train Phrase Model
Steps: (--first-step to --last-step)
(1) prepare corpus
(2) run GIZA
(3) align words
(4) learn lexical translation
(5) extract phrases
(6) score phrases
(7) learn reordering model
(8) learn generation model
(9) create decoder config file

### 3. Train the translation model (cont'd)

So, it takes a few arguments (and a few time!):

moses-scripts/training/train-model.perl -scripts-root-dir bin/moses-scripts/ -root-dir working-dir -corpus eurov4.es-en -f es -e en -alignment grow-diag-final-and -reordering msd-bidirectional-fe -lm 0:5:eurov4.en.lm:0

It generates a configuration file moses.ini needed to run the decoder where all the necessary files are specified.



### 4. Tuning of parameters with MERT

Run the Moses script mert-moses.pl (Another slow step!)

moses-scripts/training/mert-moses.pl eurov4.es-en.dev.es
eurov4.es-en.dev.en mosesdecoder/bin/moses ./model/moses.ini
--working-dir ./tuning --rootdir bin/moses-scripts/

Insert weights into configuration file with WMT10 script:

wmt10scripts/reuse-weights.perl ./tuning/moses.ini <
./model/moses.ini > moses.weight-reused.ini



- 5. Run Moses decoder on a test set
  - Tokenise and lowecase the test set as before.
  - Filter the model with Moses script.
     (mandatory for large translation tables)

moses-scripts/training/filter-model-given-input.pl ./filteredmodel
moses.weight-reused.ini testset.es

In the decoder:

mosesdecoder/bin/moses -f ./filteredmodel/moses.ini < testset.es >
testset.translated.en

## Part III

## **MT** Evaluation

### MT Evaluation basics

- Automatic Evaluation
- BLEU
- Limits of lexical similarity

















# What can achieved with automatic evaluation (as compared to manual evaluation)

- Automatic metrics notably accelerate the development cycle of MT systems:
  - Error analysis
  - System optimisation
  - System comparison

### Besides, they are

- Costless (vs. costly)
- Objective (vs. subjective)
- Reusable (vs. non-reusable)

## **Metrics based on lexical similarity** (most of the metrics!)

- Edit Distance: WER, PER, TER
- Precision: BLEU, NIST, WNM
- Recall: ROUGE, CDER
- Precision/Recall: GTM, METEOR, BLANC, SIA

## **Metrics based on lexical similarity** (most of the metrics!)

- Edit Distance: WER, PER, TER
- Precision: BLEU, NIST, WNM
- Recall: ROUGE, CDER
- Precision/Recall: GTM, METEOR, BLANC, SIA

Nowadays, BLEU is accepted as *the standard* metric.

# BLEU: a Method for Automatic Evaluation of Machine Translation

Kishore Papineni, Salim Roukos, Todd Ward, Wei-Jing Zhu IBM Research Division

"The main idea is to use a weighted average of variable length phrase matches against the reference translations. This view gives rise to a family of metrics using various weighting schemes. We have selected a promising baseline metric from this family." Candidate 1:

It is a guide to action which ensures that the military always obeys the commands of the party.

Candidate 2:

It is to insure the troops forever hearing the activity guidebook that party direct.

Candidate 1:

It is a guide to action which ensures that the military always obeys the commands of the party.

Reference 1:

It is a guide to action that ensures that the military will forever heed Party commands.

Reference 2:

It is the guiding principle which guarantees the military forces always being under the command of the Party.

Reference 3:

Candidate 1:

It is a guide to action which ensures that the military always obeys the commands of the party.

Reference 1:

It is a guide to action that ensures that the military will forever heed Party commands.

Reference 2:

It is the guiding principle which guarantees the military forces always being under the command of the Party.

Reference 3:

Candidate 2:

It is to insure the troops forever hearing the activity guidebook that party direct.

Reference 1:

It is a guide to action that ensures that the military will forever heed Party commands.

Reference 2:

It is the guiding principle which guarantees the military forces always being under the command of the Party.

Reference 3:

Precision-based measure, but:

Candidate: The the the the the the the. Reference 1: The cat is on the mat. Reference 2: There is a cat on the mat.

```
Precision-based measure, but:
```

```
Prec. =\frac{1+}{7}
```

Candidate: The the the the the the the. Reference 1: The cat is on the mat. Reference 2: There is a cat on the mat.

```
Precision-based measure, but:
```

```
Prec. =\frac{2+}{7}
```

```
Candidate:

The the the the the the the.

Reference 1:

The cat is on the mat.

Reference 2:

There is a cat on the mat.
```

```
Precision-based measure, but:
```

```
Prec. =\frac{3+}{7}
```

Candidate: The the the the the the the. Reference 1: The cat is on the mat. Reference 2: There is a cat on the mat.

```
Precision-based measure, but:
```

Prec. 
$$=\frac{4+}{7}$$

```
Candidate:

The the the the the the the.

Reference 1:

The cat is on the mat.

Reference 2:

There is a cat on the mat.
```

```
Precision-based measure, but:
```

```
Prec. =\frac{5+}{7}
```

Candidate: The the the the the the the. Reference 1: The cat is on the mat. Reference 2: There is a cat on the mat.

```
Precision-based measure, but:
```

```
Prec. =\frac{6+}{7}
```

```
Candidate:

The the the the the the the.

Reference 1:

The cat is on the mat.

Reference 2:

There is a cat on the mat.
```

```
Precision-based measure, but:
```

```
Prec. =\frac{7}{7}
```

```
Candidate:

The the the the the the the.

Reference 1:

The cat is on the mat.

Reference 2:

There is a cat on the mat.
```

A reference word should only be matched once.

Algorithm:

- Count number of times  $w_i$  occurs in each reference.
- Keep the minimum between the maximum of (1) and the number of times w<sub>i</sub> appears in the candidate (*clipping*).
- Add these values and divide by candidate's number of words.

```
Modified 1-gram precision:
```

Candidate:

The the the the the the.

Reference 1:

The cat is on the mat.

Reference 2:

There is a cat on the mat.

- $w_i \to \text{The} \\ \#_{W_i,R1} = 2 \\ \#_{W_i,R2} = 1$
- $Max_{(1)}=2, \#_{W_i,C}=7$  $\Rightarrow Min=2$
- In the second second

```
Modified 1-gram precision: P_1 =
```

Candidate: The the the the the the the. Reference 1: The cat is on the mat. Reference 2: There is a cat on the mat.

- $w_i \to \text{The} \\ \#_{W_i,R1} = 2 \\ \#_{W_i,R2} = 1$
- $Max_{(1)}=2, \#_{W_i,C}=7$  $\Rightarrow Min=2$
- In the second second

Modified 1-gram precision: 
$$P_1 = \frac{2}{2}$$

Candidate: The the the the the the the. Reference 1: The cat is on the mat. Reference 2: There is a cat on the mat.

**w**<sub>i</sub> → The #w<sub>i,R1</sub> = 2 #w<sub>i,R2</sub> = 1 **Max**<sub>(1)</sub>=2, #w<sub>i,C</sub> = 7 ⇒ Min=2
So more distinct words

Modified 1-gram precision: 
$$P_1 = \frac{2}{7}$$

Candidate: The the the the the the the. Reference 1: The cat is on the mat. Reference 2: There is a cat on the mat.

- **1**  $w_i \rightarrow \text{The}$ # $w_{i,R1} = 2$ # $w_{i,R2} = 1$ **2**  $\text{Max}_{(1)}=2, #w_i$
- $Max_{(1)} = 2, \#_{W_i,C} = 7$  $\Rightarrow Min = 2$
- On more distinct words

### Modified n-gram precision

- Straightforward generalisation to n-grams,  $P_n$ .
- Generalisation to multiple sentences:

$$P_{n} = \frac{\sum_{C \in \{\text{candidates}\}} \sum_{n \text{gram} \in C} Count_{\text{clipped}}(n \text{gram})}{\sum_{C \in \{\text{candidates}\}} \sum_{n \text{gram} \in C} Count(n \text{gram})}$$

low nhigh nadequacyfluency

### Brevity penalty

Candidate:

of the

Reference 1:

It is a guide to action that ensures that the military will forever heed Party commands.

Reference 2:

It is the guiding principle which guarantees the military forces always being under the command of the Party.

Reference 3:
# Brevity penalty

Candidate:

of the  $P_1 = 2/2, P_2 = 1/1$ 

Reference 1:

It is a guide to action that ensures that the military will forever heed Party commands.

Reference 2:

It is the guiding principle which guarantees the military forces always being under the command of the Party.

Reference 3:

It is the practical guide for the army always to heed the directions of the party.

# MT Evaluation IBM BLEU: Papineni, Roukos, Ward and Zhu [2001]

# **Brevity penalty**

$$BP = \begin{cases} 1 & \text{if } c > r \\ e^{1-r/c} & \text{if } c \le r \end{cases}$$

c candidate length, r reference length

- Multiplicative factor.
- At sentence level, huge punishment for short sentences.
- Estimated at document level.

# **BiLingual Evaluation Understudy, BLEU**

$$\mathsf{BLEU} = \mathsf{BP} \cdot \exp\left(\sum_{n=1}^{N} w_n \log P_n\right)$$

- Geometric average of  $P_n$  (empirical suggestion).
- $w_n$  positive weights summing to one.
- Brevity penalty.

# Paper's Conclusions

- BLEU correlates with human judgements.
- It can distinguish among similar systems.
- Need for multiple references or a big test with heterogeneous references.
- More parametrisation in the future.

# Watch out with BLEU implementations!

There are several widely used implementations of BLEU. (Moses multi-bleu.perl script, NIST mteval-vXX.pl script, etc.)

Results differ because of:

- Different tokenisation approach.
- Different definition of *closest reference* in the brevity penalty estimation.

**NIST** is based on BLEU but:

- Arithmetic average of *n*-gram counts rather than a geometric average.
- Informative *n*-grams are given more weight.
- Different definition of brevity penalty.



# Limits of lexical similarity

The reliability of lexical metrics depends very strongly on the heterogeneity/representativity of reference translations.

e: This sentence is going to be difficult to evaluate.

Ref1: The evaluation of the clause is complicated. Ref2: The sentence will be hard to qualify. Ref3: The translation is going to be hard to evaluate. Ref4: It will be difficult to punctuate the output.

Lexical similarity is nor a sufficient neither a necessary condition so that two sentences convey the same meaning.



# Limits of lexical similarity

The reliability of lexical metrics depends very strongly on the heterogeneity/representativity of reference translations.

e: This sentence is going to be difficult to evaluate.

Ref1: The evaluation of the clause is complicated. Ref2: The sentence will be hard to qualify. Ref3: The translation is going to be hard to evaluate. Ref4: It will be difficult to punctuate the output.

Lexical similarity is nor a sufficient neither a necessary condition so that two sentences convey the same meaning.



# Limits of lexical similarity

The reliability of lexical metrics depends very strongly on the heterogeneity/representativity of reference translations.

e: This sentence is going to be difficult to evaluate.

Ref1: The evaluation of the clause is complicated. Ref2: The sentence will be hard to qualify. Ref3: The translation is going to be hard to evaluate. Ref4: It will be difficult to punctuate the output.

Lexical similarity is nor a sufficient neither a necessary condition so that two sentences convey the same meaning.

Recent efforts to go over lexical similarity

Extend the reference material:

• Using lexical variants such as morphological variations or synonymy lookup or using paraphrasing support.

Compare other linguistic features than words:

- Syntactic similarity: shallow parsing, full parsing (constituents /dependencies).
- Semantic similarity: named entities, semantic roles, discourse representations.

Combination of the existing metrics.

Recent efforts to go over lexical similarity

Extend the reference material:

• Using lexical variants such as morphological variations or synonymy lookup or using paraphrasing support.

Compare other linguistic features than words:

- Syntactic similarity: shallow parsing, full parsing (constituents /dependencies).
- Semantic similarity: named entities, semantic roles, discourse representations.

Combination of the existing metrics.

Recent efforts to go over lexical similarity

Extend the reference material:

• Using lexical variants such as morphological variations or synonymy lookup or using paraphrasing support.

Compare other linguistic features than words:

- Syntactic similarity: shallow parsing, full parsing (constituents /dependencies).
- Semantic similarity: named entities, semantic roles, discourse representations.

Combination of the existing metrics.

### Towards Heterogeneous Automatic MT Evaluation



### Towards Heterogeneous Automatic MT Evaluation



# ASIYA

Asiya has been designed to assist both **system** and metric **developers** by offering a rich repository of metrics and meta-metrics.

http://nlp.lsi.upc.edu/asiya/

Summary

# MT Evaluation: keep in mind

- Evaluation is important in the system development cycle. Automatic evaluation accelerates significatively the process.
- Up to now, most (common) metrics rely on lexical similarity, but it cannot assure a correct evaluation.
- Current work is being devoted to go beyond lexical similarity.

# 7 MT Evaluation basics

- 8 Evaluation system
  - Software
  - Steps
  - Demo

# **Evaluate the results**

- With BLEU scoring tool. Available as a Moses script or from NIST: http://www.itl.nist.gov/iad/mig/tools/mtevalv13a-20091001.tar.gz
- With Asiya package: http://nlp.lsi.upc.edu/asiya/

# 1. Evaluate the results

With BLEU scoring tool in Moses:

moses/scripts/generic/multi-bleu.perl references.en <
testset.translated.en</pre>

### Steps

## With the Asiya toolkit:

Asiya.pl -eval single,ulc -g sys Asiya.config

#### input=raw

SRCLANG=de TRGLANG=en SRCCASE=cs TRGCASE=cs

Steps

With the Asiya toolkit:

Asiya.pl -eval single,ulc -g sys Asiya.config

### System evaluation with Asiya

Asiya.pl -eval single,ulc -m metrSet Asiya.config

SRCLANG=de TRGLANG=en

metrSet=1-PER 1-TER 1-WER BLEU-4 CP-Oc-\* CP-Op-\* CP-STM-9 DP-HWC-c-4 DP-HWC-r-4 DP-HWC-w-4 DP-Oc-\* DP-Ol-\* DP-Dr-\* DR-Or-\* DR-Orp-\* DR-STM-9 GTM-1 GTM-2 GTM-3 MTR-exact MTR-wisem MTR-wisem MTR-wisem MTR-wisem NE-Me-\* NE-Oe-\* NE-Oe-\*\* NIST-5 RG-L RG-S\* RG-SU\* RG-W1-2 SP-Oc-\* SP-Op-\* SP-cNIST-5 SP-iobNIST-5 SP-INIST-5 SP-DNIST-5 SR-Mr-\* SR-Mrv\* SR-Or SR-Or-\* SR-Orv

### Metrics in Asiya (English)

#### METRIC NAMES

------

668 metrics are available for language 'en

METRICS = { -PER, -TER, -TERbase, -TERb, -TERbase, -TERb, -TERbase, -TERb, -WER, BLEU, BLEU-1, BLEU-2, BLEU-3, BLEU-4, BLEU-2, BLEU-3, BLEU-4, CP-0c(\*), CP-0c(\*), CP-0c(ADJP), CP-0c(CONJP), CP-0c(FRA G), CP-Oc(INTJ), CP-Oc(LST), CP-Oc(NAC), CP-Oc(NP), CP-Oc(NX), CP-Oc(O), CP-Oc(PP), CP-Oc(PRI), CP-Oc(PRI), CP-Oc(OP), CP-Oc(SI), CP-Oc(SINV), CP-OC -Oc(UCP), CP-Oc(VP), CP-Oc(WHADJP), CP-Oc(WHADJP), CP-Oc(WHAPP), CP-Oc(X), CP-Op(#), CP-Op(\*), C -Op(CC), CP-Op(CD), CP-Op(DT), CP-Op(F), CP-Op(F), CP-Op(F), CP-Op(IN), CP-Op(J), CP-Op(JJ), CP-Op(JJ), CP-Op(JJS), CP-Op(MD), CP-Op(MD), CP-Op(NN), CP-Op(NNP), C NNPS), CP-0p(NNS), CP-0p(PDT), CP-0p(PDT), CP-0p(PDS), CP-0p(PRPs), CP-0p(PRP), CP-0p(RB), CP-0p(RB), CP-0p(RBS), CP-0p(RB), CP-0p(SYM), CP-0p(T0), CP-0p(T0), CP-0p(T0), CP-0p(V), CP-0p( -Op(VB), CP-Op(VBD), CP-Op(VBC), CP-Op(VBC), CP-Op(VBC), CP-Op(VBC), CP-Op(VBC), CP-Op(WDC), CP-Op(WPC), CP-Op(WPC), CP-Op(WRC), CP-Op(VC), CP-STN-5, CP-STM-6, CP-STM-7, CP-STM-8, CP-STM-9, CP-STM1-2, CP-STM1-3, CP-STM1-4, CP-STM1-5, CP-STM1-6, CP-STM1-7, CP-STM1-9, DP-HWCM c-1, DP-HWCM c-2, DP-HWCM c-3, DP-HWC M c-4, DP-HWCM r-1, DP-HWCM r-2, DP-HWCM r-3, DP-HWCM r-4, DP-HWCM w-1, DP-HWCM w-2, DP-HWCM w-4, DP-HWCMi c-2, DP-HWCMi c-3, DP-HWCMi c-4, DP-HWCMi r-2, DP-HWCMi r-3, DP-HWCM1\_r-4, DP-HWCM1\_w-2, DP-HWCM1\_w-3, DP-HWCM1\_w-4, DP-Oc(\*), DP-Oc(a), DP-Oc(ax), DP-Oc(aux), DP-Oc(be), DP-Oc(c), DP-Oc(comp), DP-Oc(det), DP-Oc(have), DP-Oc(n), DP-Oc(postd et), DP-Oc(pospec), DP-Oc(yredet), DP-Oc(serida), DP-Oc(serida), DP-Oc(serida), DP-Oc(subi), DP-Oc(that), DP-Oc(y), DP-Oc(ybe), DP-Oc(xsaid), DP-Ol(\*), DP-Ol(1), DP-Ol(\*), DP-O 2), DP-01(3), DP-01(4), DP-01(5), DP-01(6), DP-01(7), DP-01(8), DP-01(9), DP-0r(\*), DP-0r(anod), DP-0r(anount-value), DP-0r(appo), DP-0r(appo-mod), DP-0r(as-arg), DP-0r(as1), DP-0r(\*), D (as2), DP-Or(aux), DP-Or(be), DP-Or(being), DP-Or(by-subi), DP-Or(c), DP-Or(cn), DP-Or(compl), DP-Or(compl), DP-Or(desc), DP-Or(dest), DP-Or(det), DP-Or(else), DP-Or(fc), DP-Or , DP-Or(guest), DP-Or(have), DP-Or(head), DP-Or(i), DP-Or(inv-aux), DP-Or(inv-have), DP-Or(lex-dep), DP-Or(lex-mod), DP-Or(mod-before), DP-Or(neg), DP-Or(nn), DP-Or(num ), DP-Or(num-mod), DP-Or(obi), DP-Or(obi), DP-Or(p), DP-Or(p-spec), DP-Or(pcomp-c), DP-Or(pcomp-n), DP-Or(person), DP-Or(poss), DP-Or(post), DP-Or(post), DP-Or(print), DP-Or(post), DP-Or( DP-Or(pred), DP-Or(punc), DP-Or(rel), DP-Or(s), DP-Or(sc), DP-Or(subcat), DP-Or(subclass), DP-Or(subj), DP-Or(title), DP-Or(vrel), DP-Or(wha), DP-Or(w , DPm-HWCM c-2, DPm-HWCM c-3, DPm-HWCM c-4, DPm-HWCM r-1, DPm-HWCM r-2, DPm-HWCM r-3, DPm-HWCM r-4, DPm-HWCM w-1, DPm-HWCM w-3, DPm-HWCM w-4, DPm-HWCM i c-2, DPm-HWCM r-4, DPM-HWCM r-4 c-3, DPm-HWCMi c-4, DPm-HWCMi r-2, DPm-HWCMi r-3, DPm-HWCMi r-4, DPm-HWCMi w-2, DPm-HWCMi w-3, DPm-HWCMi w-4, DPm-Oc(\*), DPm-Oc(\*), DPm-Ol(\*), DPm-Ol(\*), DPm-Ol(1), DPm-Ol(2), DPm-Ol(3) , DPm-01(4), DPm-01(5), DPm-01(6), DPm-01(7), DPm-01(8), DPm-01(9), DPm-0r(.....), DR-Fr(\*), DR-Fr(\*), DR-0r(\*), DR-0r(\*), DR-0r(\*) b, DR-0r(\*) i, DR-0r(alfa), DR-0r(car d), DR-Or(drs), DR-Or(eq), DR-Or(imp), DR-Or(merge), DR-Or(named), DR-Or(not), DR-Or(or), DR-Or(ored), DR-Or(prop), DR-Or(rel), DR-Or(smerge), DR-Or(timex), DR-Or(who), DR-Or(or), DR-Or(or), DR-Or(erge), DR-Or(smerge), DR-Or(smerge DR-Orp(\*), DR-Orp(\*) b, DR-Orp(\*) i, DR-Orp(alfa), DR-Orp(card), DR-Orp(dr), DR-Orp(drs), DR-Orp(eq), DR-Orp(imp), DR-Orp(merge), DR-Orp(mamed), DR-Orp(not), DR-Orp(or), DR-Orp(or) ed), DR-Orp(prop), DR-Orp(rel), DR-Orp(smerge), DR-Orp(timex), DR-Orp(who), DR-Pr(\*), DR-Prp(\*), DR-Rrp(\*), DR-SrM-1, DR-STM-2, DR-STM-3, DR-STM-4, DR-STM-4 b, DR-STM-4 i , DR-STM-5, DR-STM-6, DR-STM-7, DR-STM-8, DR-STM-9, DR-STM-2, DR-STM1-3, DR-STM1-4, DR-STM1-5, DR-STM1-6, DR-STM1-7, DR-STM1-8, DR-STM1-9, DRdoc-0r(\*), DRdoc-0r(\*) b, DR doc-Or(\*) i, DRdoc-Or(alfa), DRdoc-Or(card), DRdoc-Or(dr), DRdoc-Or(drs), DRdoc-Or(eg), DRdoc-Or(imp), DRdoc-Or(merge), DRdoc-Or(named), DRdoc-Or(not), DRdoc-Or(or), DRdoc-Or(ored) , DRdoc.Or(prop), DRdoc-Or(rel), DRdoc-Or(smerge), DRdoc-Or(timex), DRdoc-Or(who), DRdoc-Orp(\*), DRdoc-Orp(\*) b, DRdoc-Orp(\*) i, DRdoc-Orp(alfa), DRdoc-Orp(card), DRdoc-Orp(dr), DR doc-Orp(drs), DRdoc-Orp(eq), DRdoc-Orp(imp), DRdoc-Orp(merge), DRdoc-Orp(named), DRdoc-Orp(not), DRdoc-Orp(or), DRdoc-Orp(pred), DRdoc-Orp(pred), DRdoc-Orp(rel), DRdoc-Orp(merge), DRdoc-Orp(timex), DRdoc-Orp(wha), DRdoc-STM-1, DRdoc-STM-2, DRdoc-STM-3, DRdoc-STM-4, DRdoc-STM-4 b, DRdoc-STM-4 i, DRdoc-STM-5, DRdoc-STM-6, DRdoc-STM-7, DRdoc-STM-8, DRdoc-STM-9, DRdoc-, DRdoc-STMi-2, DRdoc-STMi-3, DRdoc-STMi-4, DRdoc-STMi-5, DRdoc-STMi-6, DRdoc-STMi-7, DRdoc-STMi-8, DRdoc-STMi-9, FL, GTM-1, GTM-2, GTM-3, METEOR-ex, METEOR-ex, METEOR-st, METE v, NE-Me(\*), NE-Me(ANGLE QUANTITY), NE-Me(DATE), NE-Me(DISTANCE QUANTITY), NE-Me(LANGUAGE), NE-Me(LOC), NE-Me(MEASURE), NE-Me(METHOD), NE-Me(MONEY), NE-Me(NUM), NE-ME ORG), NE-Me(PER), NE-Me(PERCENT), NE-Me(PROJECT), NE-Me(SIZE QUANTITY), NE-Me(SPEED QUANTITY), NE-Me(SYSTEM), NE-Me(TEMPERATURE QUANTITY), NE-Me(TIME), NE-Me(WEIGHT QUANTITY), NE-O e(\*), NE-Oe(\*\*), NE-Oe(ANGLE QUANTITY), NE-Oe(DATE), NE-Oe(DISTANCE QUANTITY), NE-Oe(LANGUAGE), NE-Oe(NOC), NE-Oe(MERSURE), NE-Oe(MISC), NE-Oe(MONEY), NE-Oe(NONEY), NE-OE -Oe(O), NE-Oe(ORG), NE-Oe(PER), NE-Oe(PERCENT), NE-Oe(PROJECT), NE-Oe(SIZE OUANTITY), NE-Oe(SPEED OUANTITY), NE-Oe(SYSTEM), NE-Oe(TEMPERATURE OUANTITY), NE-OE(TE UANTITY), NIST, NIST-1, NIST-2, NIST-3, NIST-4, NIST-5, NIST1-2, NIST1-3, NIST1-4, NIST1-5, 01, P1, ROUGE-1, ROUGE-2, ROUGE-4, ROUGE-4, ROUGE-5, RO P-0c(\*), SP-0c(ADJP), SP-0c(ADJP), SP-0c(CONJP), SP-0c(INTJ), SP-0c(INT), SP-0c(NP), SP-0c(OP), SP-0c(PP), SP-0c(SBAR), SP-0c(VP), SP-0c(VP), SP-0c(VP), SP-0c(SBAR), SP-0c(VP), SP-0c(VP), SP-0c(SBAR), ''), SP-0p((), SP-0p()), SP-0p(\*), SP-0p(,), SP-0p(,), SP-0p(2), SP-0p(CD), SP-0p(CD), SP-0p(EX), SP-0p(FN), SP-0p(FN), SP-0p(JN), SP-0p(JJ), S JJS), SP-0p(LS), SP-0p(ND), SP-0p(N), SP-0p(NN), SP-0p(NNP), SP-0p(NNPS), SP-0p(NNS), SP-0p(P), SP-0p(PDT), SP-0p(PRPS), SP-0p(PRPS), SP-0p(PRP), SP-0p(RB), SP-0p(RB -Op(RBS), SP-Op(RP), SP-Op(VBZ), SP-Op(VD), SP-Op(WD), SP-Op(VB), SP-Op(VBD), SP-Op(VBD), SP-Op(VBD), SP-Op(VBP), SP-Op(VBZ), SP-Op(WD), SP-Op(WDT), S SP-OD(WRB), SP-OD(\*), SP-CNIST, SP-CNIST-1, SP-CNIST-2, SP-CNIST-3, SP-CNIST-4, SP-CNIST-5, SP-CNIST-3, SP-CNIST-4, SP-CNIST-4 ST-2, SP-iobNIST-3, SP-iobNIST-4, SP-iobNIST-5, SP-iobNISTi-2, SP-iobNISTi-3, SP-iobNISTi-4, SP-iobNISTi-5, SP-UNIST-3, SP-UNIST-1, SP-UNIST-3, SP-UNIST-3, SP-UNIST-5, SP--UNISTI-2, SP-UNISTI-3, SP-UNISTI-4, SP-UNISTI-5, SP-DNIST-1, SP-DNIST-2, SP-DNIST-3, SP-DNIST-5, SP-DNIST-5, SP-DNIST-3, SP-DNISTI-5, SP-DNIST-5, SP-, SR-MFr(\*), SR-MPr(\*), SR-Mr(\*), SR-Mr(\*), SR-Mr(\*) b, SR-Mr(\*) 1, SR-Mr(A0), SR-Mr(A1), SR-Mr(A2), SR-Mr(A3), SR-Mr(A3), SR-Mr(A5), SR-Mr(AA), SR-Mr(AA) r(AM-DIR), SR-Mr(AM-DIS), SR-Mr(AM-EXT), SR-Mr(AM-LOC), SR-Mr(AM-MNR), SR-Mr(AM-MOD), SR-Mr(AM-NEG), SR-Mr(AM-PNC), SR-Mr(AM-REC), SR-Mr(AM-TMP), SR-Mrv(\*), SR-Mrv(\* ) b, SR-Mrv(\*) 1, SR-Mrv(A0), SR-Mrv(A1), SR-Mrv(A2), SR-Mrv(A3), SR-Mrv(A4), SR-Mrv(A5), SR-Mrv(AA), SR-Mrv(AM-ADV), SR-Mrv(AM-CAU), SR-Mrv(AM-DIS), SR-Mrv(AM-EXT) SR-Mrv(AM-LOC), SR-Mrv(AM-NNR), SR-Nrv(AM-NOD), SR-Mrv(AM-NEG), SR-Mrv(AM-PNC), SR-Mrv(AM-PRD), SR-Mrv(AM-REC), SR-Mrv(AM-TMP), SR-Nv, SR-O1, SR-Or(+), SR-Or(+), SR-Or(+), SR-Or(+), SR-O1, SR ) 1, SR-0r(A0), SR-0r(A1), SR-0r(A2), SR-0r(A3), SR-0r(A4), SR-0r(A5), SR-0r(AA), SR-0r(AM-ADV), SR-0r(AM-ADV), SR-0r(AM-DIS), SR-0r(AM-DIS), SR-0r(AM-EXT), SR-0r(AM-LOC), SR-0r(AM-ADV), SR-0r(AM-ADV), SR-0r(AM-DIS), SR-0r(AM-DIS), SR-0r(AM-ADV), SR-0r(ADV), SR-0r -INR), SR-Or(AM-MOD), SR-Or(AM-NEG), SR-Or(AM-PRC), SR-Or(AM-PRD), SR-Or(AM-REC), SR-Or(AM-TMP), SR-Or 1, SR-Or 1, SR-Or 1, SR-Or (\*), SR-OR (\* 1), SR-0rv(A2), SR-0rv(A3), SR-0rv(A4), SR-0rv(A5), SR-0rv(AA), SR-0rv(AM-ADV), SR-0rv(AM-CAU), SR-0rv(AM-DIR), SR-0rv(AM-DIS), SR-0rv(AM-EXT), SR-0rv(AM-LOC), SR-0rv(AM-NNR), SR-0rv(AM-NNR), SR-0rv(AM-DIS), SR-0rv(AM-DIS) ry(AM-MOD), SR-Ory(AM-NEG), SR-Ory(AM-PNC), SR-Ory(AM-PRD), SR-Ory(AM-REC), SR-Ory(AM-TMP), SR-Ory b, SR-Ory 1, SR-O

# 2. Evaluate the results on-line

Asiya Interface

http://asiya.lsi.upc.edu/demo/asiya\_online.php

# 3. Analise the results on-line

t-Search Interface

http://asiya.lsi.upc.edu/demo/tsearch\_upload.php

### Demo: http://asiya.lsi.upc.edu/demo/asiya\_online.php

🥪 🗇 🛛 Asiya: An Open Toolkit for Automatic Machine Translation (Meta-)Evaluation - Mozilla Firefox				
<u>F</u> itxer <u>E</u> dita <u>V</u> isualitza Hi <u>s</u> torial Ad <u>r</u> eces d'interès Ei <u>n</u> es Ajuda				
🛄 Asiya: An Open Toolkit for Aut 👍				
🔶 🕑 asiya.lsi.upc.edu/demo/asiya_online.php 🔅 🕫 🚷 🛪 Google			🧟 👆 🔕	
<b>Asiya -</b> An Online	<b>Online</b> Toolkit for Automatic Machine Translation i	Evaluation		Â
Asiya Tes	bed Data:	Asiya Files	Edit View Tools Help	
Data Format Input format: Input already to	raw 🗘 Source Language: other 🗘	Source Case: case sensitive \$		=
Files Source file: Source text:	Navega No sha seleccionat cap fitw Write some text here instead of uplos a file.	er. Upload		
Nafeence file: Navega No sha seleccional cap fiber. Upload Reference ten: File. No sha seleccional cap fiber.				
Translation Sys	em Mes: Navega) No sha seleccional cap Mixe em text Write some text here instead of uplos a file.	.di Upload iding		

An experiment: http://asiya.lsi.upc.edu/demo/userexperiment.php

# Was it easy?

Prepared experiment

http://asiya.lsi.upc.edu/demo/userexperiment.php

# Part IV

# Appendix: References

# History of SMT

- Weaver, 1949 [Wea55]
- Alpac Memorandum [Aut66]
- Hutchins, 1978 [Hut78]
- Slocum, 1985 [Slo85]

# The beginnings, word-based SMT

- Brown et al., 1990 [BCP+90]
- Brown et al., 1993 [BPPM93]

# Phrase-based model

- Och et al., 1999 [OTN99]
- Koehn et al, 2003 [KOM03]

# Log-linear model

- Och & Ney, 2002 [ON02]
- Och & Ney, 2004 [ON04]

# Factored model

• Koehn & Hoang, 2007 [KH07]

# Syntax-based models

- Yamada & Knight, 2001 [YK01]
- Chiang, 2005 [Chi05]
- Carreras & Collins, 2009 [CC09]

# **Discriminative models**

- Carpuat & Wu, 2007 [CW07]
- Bangalore et al., 2007 [BHK07]
- Giménez & Màrquez, 2008 [GM08]

# Language model

• Kneser & Ney, 1995 [KN95]

# MERT

• Och, 2003 [Och03]

# **Domain adaptation**

• Bertoldi and Federico, 2009 [Och03]

# Reordering

- Crego & Mariño, 2006 [Cn06]
- Bach et al., 2009 [BGV09]
- Chen et al., 2009 [CWC09]

# Systems combination

- Du et al., 2009 [DMW09]
- Li et al., 2009 [LDZ+09]
- Hildebrand & Vogel, 2009 [HV09]

# Alternative systems in development

- Blunsom et al., 2008 [BCO08]
- Canisius & van den Bosch, 2009 [CvdB09]
- Chiang et al., 2009 [CKW09]
- Finch & Sumita, 2009 [FS09]
- Hassan et al., 2009 [HSW09]
- Shen et al., 2009 [SXZ+09]

# Evaluation

- Papineni, 2002 [PRWZ02]
- Doddington, 2002 [Dod02]
- Banerjee & Alon Lavie, 2005 [BL05]
- Giménez & Amigó, 2006 [GA06]

# Surveys, theses and tutorials

• Knight, 1999

http://www.isi.edu/natural-language/mt/wkbk.rtf

• Knight & Koehn, 2003

http://people.csail.mit.edu/people/koehn/publications/tutorial2003.pdf

• Koehn, 2006

 $http://www.iccs.informatics.ed.ac.uk/\ pkoehn/publications/tutorial2006.pdf$ 

• Way & Hassan, 2009

http://www.medar.info/conference\_all/2009/Tutorial\_3.pdf

- Lopez, 2008 [Lop08]
- Giménez, 2009 [Gim08]
# References I

Automatic Language Processing Advisory Committee (ALPAC). Language and Machines. Computers in Translation and Linguistics. Technical Report Publication 1416, Division of Behavioural Sciences, National Academy of Sciences, National Research Council, Washington, D.C., 1966.

Phil Blunsom, Trevor Cohn, and Miles Osborne. A discriminative latent variable model for statistical machine translation. In ACL-08: HLT. 46th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies, pages 200–208, 2008.



Peter F. Brown, John Cocke, Stephen A. Della Pietra, Vincent J. Della Pietra, Fredrick Jelinek, John D. Lafferty, Robert L. Mercer, and Paul S. Roossin. A statistical approach to machine translation. *Computational Linguistics*, 16(2):79–85, 1990.

Nguyen Bach, Qin Gao, and Stephan Vogel.

Source-side dependency tree reordering models with subtree movements and constraints.

In Proceedings of the Twelfth Machine Translation Summit (MTSummit-XII), Ottawa, Canada, August 2009. International Association for Machine Translation.

# **References II**

### Srinivas Bangalore, Patrick Haffner, and Stephan Kanthak.

Statistical Machine Translation through Global Lexical Selection and Sentence Reconstruction.

In Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics (ACL), pages 152–159, 2007.

### Satanjeev Banerjee and Alon Lavie.

METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments.

In Proceedings of ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for MT and/or Summarization, 2005.



Peter F. Brown, Vincent J. Della Pietra, Stephen A. Della Pietra, and Robert L. Mercer.

The mathematics of statistical machine translation: parameter estimation. *Computational Linguistics*, 19(2):263–311, 1993.

### Xavier Carreras and Michael Collins.

Non-projective parsing for statistical machine translation. In Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 200–209, Singapore, August 2009.

# **References III**



## David Chiang.

A hierarchical phrase-based model for statistical machine translation. In Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL'05), pages 263–270, Ann Arbor, Michigan, June 2005. Association for Computational Linguistics.

David Chiang, Kevin Knight, and Wei Wang. 11,001 new features for statistical machine translation.

In NAACL '09: Human Language Technologies: the 2009 annual conference of the North American Chapter of the ACL, pages 218–226. Association for Computational Linguistics, 2009.



Josep M<sup>a</sup> Crego and José B. Mari no. Improving smt by coupling reordering and decoding. *Machine Translation*, 20(3):199–215, March 2006.

Sander Canisius and Antal van den Bosch. A constraint satisfaction approach to machine translation.

In Lluís Màrquez and Harold Somers, editors, *EAMT-2009: Proceedings of the* 13th Annual Conference of the European Association for Machine Translation, pages 182–189, 2009.

# **References IV**

## 

### Marine Carpuat and Dekai Wu.

Improving Statistical Machine Translation Using Word Sense Disambiguation. In Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 61–72, 2007.



Han-Bin Chen, Jian-Cheng Wu, and Jason S. Chang. Learning bilingual linguistic reordering model for statistical machine translation. In NAACL '09: Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics, pages 254–262, Morristown, NJ, USA, 2009. Association for Computational Linguistics.



## Jinhua Du, Yanjun Ma, and Andy Way.

Source-side context-informed hypothesis alignment for combining outputs from Machine Translation systems.

In Proceedings of the Machine Translation Summit XII, pages 230–237, Ottawa, ON, Canada., 2009.



### George Doddington.

Automatic evaluation of machine translation quality using n-gram co-occurrence statistics.

In Proceedings of the 2nd Internation Conference on Human Language Technology, pages 138–145, 2002.

# References V

### Andrew Finch and Eiichiro Sumita.

Bidirectional phrase-based statistical machine translation.

In Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing, pages 1124–1132, Singapore, August 2009. Association for Computational Linguistics.



Jesús Giménez and Enrique Amigó. IQMT: A Framework for Automatic Machine Translation Evaluation. In *Proceedings of the 5th LREC*, pages 685–690, 2006.



### Jesś Giménez.

*Empirical Machine Translation and its Evaluation.* PhD thesis, Universitat Politècnica de Catalunya, July 2008.



Jesús Giménez and Lluís Màrquez. *Discriminative Phrase Selection for SMT*, pages 205–236. NIPS Workshop Series. MIT Press, 2008.

Hany Hassan, Khalil Sima'an, and Andy Way. A syntactified direct translation model with linear-time decoding. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing*, pages 1182–1191, Singapore, August 2009. Association for Computational Linguistics.

# **References VI**



## W. J. Hutchins.

Machine translation and machine-aided translation. *Journal of Documentation*, 34(2):119–159, 1978.



# Almut Silja Hildebrand and Stephan Vogel. CMU system combination for WMT'09.

In Proceedings of the Fourth Workshop on Statistical Machine Translation, pages 47–50, Athens, Greece, March 2009. Association for Computational Linguistics.

### Philipp Koehn and Hieu Hoang. Factored Translation Models.

In Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 868–876, 2007.



### R. Kneser and H. Ney.

Improved backing-off for m-gram language modeling. *icassp*, 1:181–184, 1995.

### Philipp Koehn, Franz Josef Och, and Daniel Marcu. Statistical phrase-based translation.

In Proceedings of the Human Language Technology and North American Association for Computational Linguistics Conference (HLT/NAACL), Edomonton, Canada, May 27-June 1 2003.

# **References VII**

Mu Li, Nan Duan, Dongdong Zhang, Chi-Ho Li, and Ming Zhou.

Collaborative decoding: Partial hypothesis re-ranking using translation consensus between decoders.

In Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP, pages 585–592, Suntec, Singapore, August 2009. Association for Computational Linguistics.

10	
1.5	_

### Adam Lopez.

Statistical machine translation. ACM Comput. Surv., 40(3), 2008.



### Franz Josef Och.

Minimum error rate training in statistical machine translation.

In Proc. of the Association for Computational Linguistics, Sapporo, Japan, July 6-7 2003.



### Franz Josef Och and Hermann Ney.

Discriminative Training and Maximum Entropy Models for Statistical Machine Translation.

In Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics (ACL), pages 295–302, 2002.

### Franz Josef Och and Hermann Ney.

The alignment template approach to statistical machine translation. *Computational Linguistics*, 30(4):417–449, 2004.



Franz Josef Och, Christoph Tillmann, and Hermann Ney. Improved alignment models for statistical machine translation. In Proc. of the Conference on Empirical Methods in Natural Language Processing and Very Large Corpora, pages 20–28, University of Maryland, College Park, MD, June 1999.



Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the Association of Computational Linguistics*, pages 311–318, 2002.



### Jonathan Slocum.

A survey of machine translation: its history, current status, and future prospects. Comput. Linguist., 11(1):1–17, 1985.

Libin Shen, Jinxi Xu, Bing Zhang, Spyros Matsoukas, and Ralph Weischedel. Effective use of linguistic and contextual information for statistical machine translation.

In Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing, pages 72–80, Singapore, August 2009. Association for Computational Linguistics.



### Warren Weaver.

### Translation.

In William N. Locke and A. Donald Boothe, editors, *Machine Translation of Languages*, pages 15–23. MIT Press, Cambridge, MA, 1949/1955. Reprinted from a memorandum written by Weaver in 1949.

### Kenji Yamada and Kevin Knight.

### A syntax-based statistical translation model.

In Proceedings of the 39rd Annual Meeting of the Association for Computational Linguistics (ACL'01), Toulouse, France, July 2001.