

Statistical Machine Translation

A practical tutorial

Cristina España i Bonet
LSI Department
Universitat Politècnica de Catalunya

MOLTO Kickoff Meeting
UPC, Barcelona
11th March, 2010

- 1 Introduction
- 2 Basics
- 3 Components
- 4 The log-linear model
- 5 Beyond standard SMT
- 6 MT Evaluation

Part I: SMT background

~ 90 minutes

7 Translation system

Part II: SMT experiments

8 Evaluation system

~ 30 minutes

9 References

Part III: References

Part I

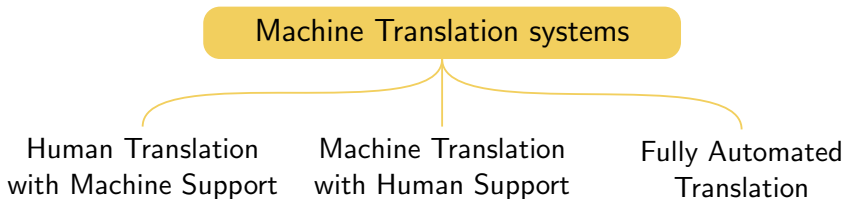
SMT background

Outline

- 1 Introduction
- 2 Basics
- 3 Components
- 4 The log-linear model
- 5 Beyond standard SMT
- 6 MT Evaluation

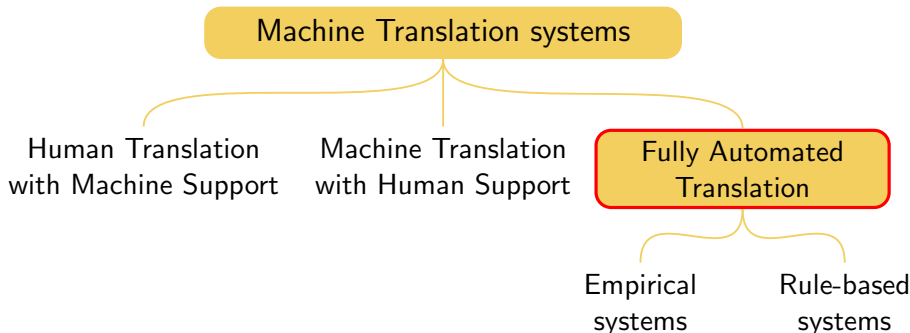
Introduction

Machine Translation Taxonomy



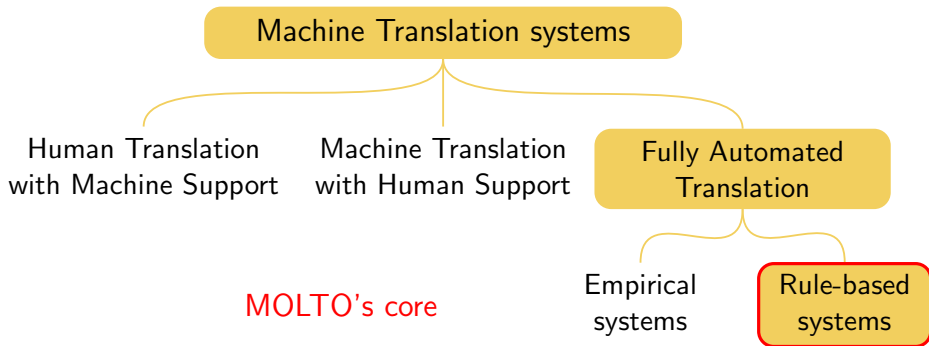
Introduction

Machine Translation Taxonomy



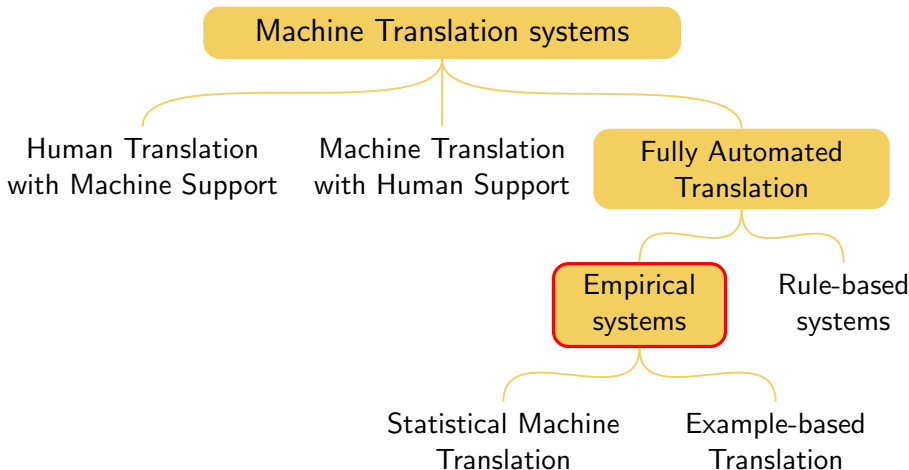
Introduction

Machine Translation Taxonomy



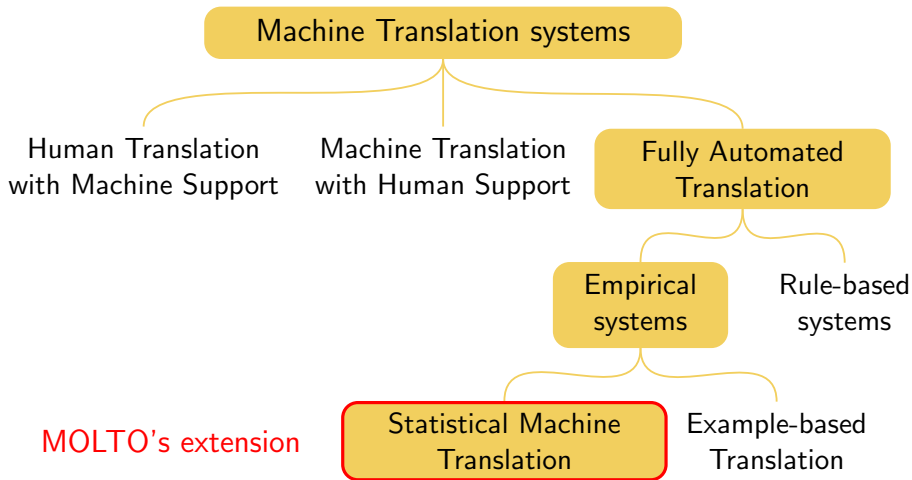
Introduction

Machine Translation Taxonomy



Introduction

Machine Translation Taxonomy



Introduction

Empirical Machine Translation

Empirical MT relies on large parallel aligned corpora.

L'objectiu de MOLTO és desenvolupar un conjunt d'eines per a traduir textos entre diversos idiomes en temps real i amb alta qualitat. Les llengües són mòduls separats en l'eina i per tant es poden canviar; els prototips que es construïran cobriran la major part dels 23 idiomes oficials de la UE.

Com a tècnica principal, MOLTO utilitza gramàtiques semàntiques de domini específic i interlingües basades en ontologies. Aquests components s'implementen en GF (Grammatical Framework), un formalisme de gramàtiques on es relacionen diversos idiomes a través d'una sintaxi abstracta comú. El GF s'ha aplicat en diversos dominis de mida petita i mitjana, típicament per tractar fins a un total de deu idiomes, però MOLTO ampliarà això en termes de productivitat i aplicabilitat.

Part de l'ampliació es dedicarà a augmentar la mida dels dominis i el nombre d'idiomes. Una part important és fer la tecnologia accessible per als experts del domini sense experiència amb GFs i reduir al mínim l'esforç necessari per a la construcció d'un traductor. Idealment, això es pot fer només estenent un llexicó i escrivint un conjunt de frases d'exemple.

MOLTO's goal is to develop a set of tools for translating texts between multiple languages in real time with high quality. Languages are separate modules in the tool and can be varied; prototypes covering a majority of the EU's 23 official languages will be built.

As its main technique, MOLTO uses domain-specific semantic grammars and ontology-based interlinguas. These components are implemented in GF (Grammatical Framework), which is a grammar formalism where multiple languages are related by a common abstract syntax. GF has been applied in several small-to-medium size domains, typically targeting up to ten languages but MOLTO will scale this up in terms of productivity and applicability.

A part of the scale-up is to increase the size of domains and the number of languages. A more substantial part is to make the technology accessible for domain experts without GF expertise and minimize the effort needed for building a translator. Ideally, this can be done by just extending a lexicon and writing a set of example sentences.

Introduction

Empirical Machine Translation

Empirical MT relies on large parallel aligned corpora.

L'objectiu de MOLTO és desenvolupar un conjunt d'eines per a traduir textos entre diversos idiomes en temps real i amb alta qualitat. Les llengües són mòduls separats en l'eina i per tant es poden canviar; els prototips que es construïran cobriran la major part dels 23 idiomes oficials de la UE.

Com a tècnica principal, MOLTO utilitza gramàtiques semàntiques de domini específic i interlingües basades en ontologies. Aquests components s'implementen en GF (Grammatical Framework), un formalisme de gramàtiques on es relacionen diversos idiomes a través d'una sintaxi abstracta comú. El GF s'ha aplicat en diversos dominis de mida petita i mitjana, típicament per tractar fins a un total de deu idiomes, però MOLTO ampliarà això en termes de productivitat i aplicabilitat.

Part de l'ampliació es dedicarà a augmentar la mida dels dominis i el nombre d'idiomes. Una part important és fer la tecnologia accessible per als experts del domini sense experiència amb GFs i reduir al mínim l'esforç necessari per a la construcció d'un traductor. Idealment, això es pot fer només estenent un llexicó i escrivint un conjunt de frases d'exemple.

MOLTO's goal is to develop a set of tools for translating texts between multiple languages in real time with high quality. languages are separate modules in the tool and can be varied; prototypes covering a majority of the EU's 23 official languages will be built.

As its main technique, MOLTO uses domain-specific semantic grammars and ontology-based interlinguas. These components are implemented in GF (Grammatical Framework), which is a grammar formalism where multiple languages are related by a common abstract syntax. GF has been applied in several small-to-medium size domains, typically targeting up to ten languages but MOLTO will scale this up in terms of productivity and applicability.

A part of the scale-up is to increase the size of domains and the number of languages. A more substantial part is to make the technology accessible for domain experts without GF expertise and minimize the effort needed for building a translator. Ideally, this can be done by just extending a lexicon and writing a set of example sentences.

Introduction

Empirical Machine Translation

Aligned parallel corpora numbers

Corpora

Corpus	# segments (app.)	# words (app.)
JRC-Acquis	$1.0 \cdot 10^6$	$30 \cdot 10^6$
Europarl	$1.5 \cdot 10^6$	$45 \cdot 10^6$
United Nations	$3.8 \cdot 10^6$	$100 \cdot 10^6$

Books

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

Introduction

Empirical Machine Translation

Aligned parallel corpora numbers

Corpora

Corpus	# segments (app.)	# words (app.)
JRC-Acquis	$1.0 \cdot 10^6$	$30 \cdot 10^6$
Europarl	$1.5 \cdot 10^6$	$45 \cdot 10^6$
United Nations	$3.8 \cdot 10^6$	$100 \cdot 10^6$

Books

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

Outline

- 1 Introduction
- 2 Basics**
- 3 Components
- 4 The log-linear model
- 5 Beyond standard SMT
- 6 MT Evaluation

SMT, basics

The beginnings, summarised timeline

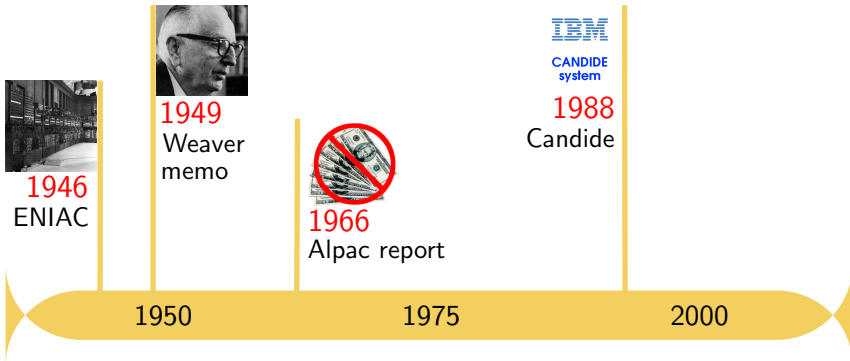


A horizontal timeline bar with a yellow-to-orange gradient, featuring arrowheads at both ends. It is marked with the years 1950, 1975, and 2000. A vertical line extends upwards from the timeline at a point between 1975 and 2000, with the year 1988 written in red above it. To the right of the vertical line, the text 'IBM CANDIDE system' is displayed in blue, with 'IBM' in a larger, bold font. Below this, the word 'Candide' is written in black.

IBM
CANDIDE
system
1988
Candide

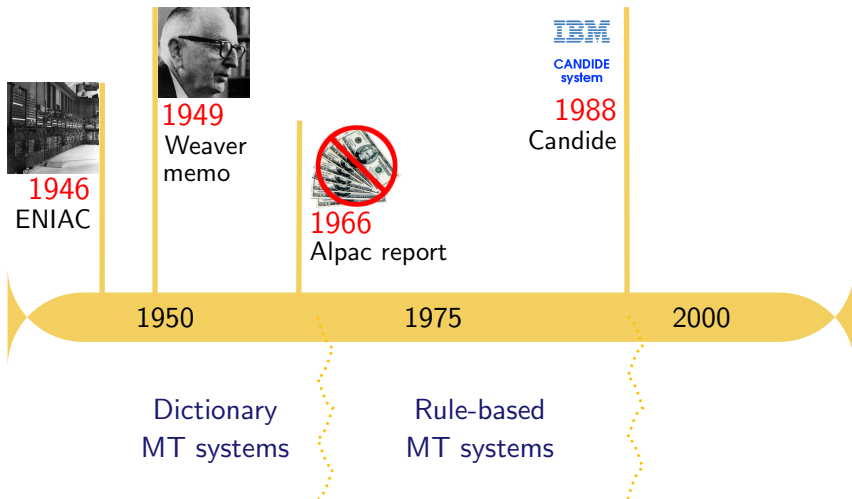
SMT, basics

The beginnings, summarised timeline



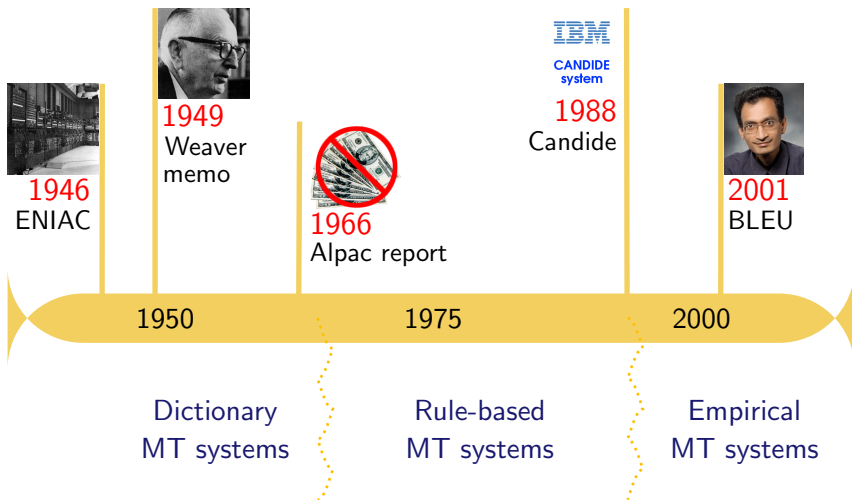
SMT, basics

The beginnings, summarised timeline



SMT, basics

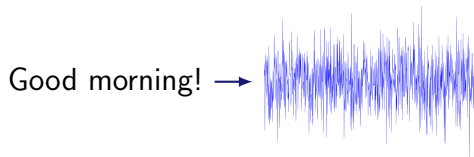
The beginnings, summarised timeline



SMT, basics

The Noisy Channel approach

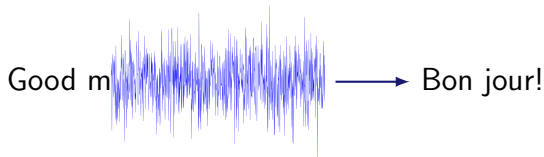
The Noisy Channel as a statistical approach to translation:



SMT, basics

The Noisy Channel approach

The Noisy Channel as a statistical approach to translation:



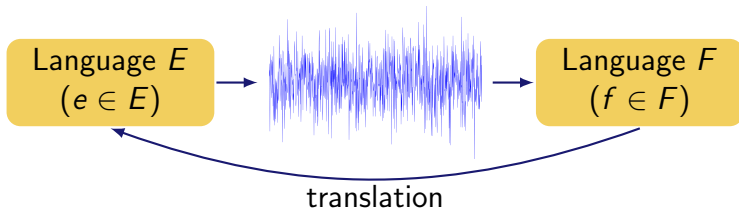
SMT, basics

The Noisy Channel approach

The Noisy Channel as a statistical approach to translation:

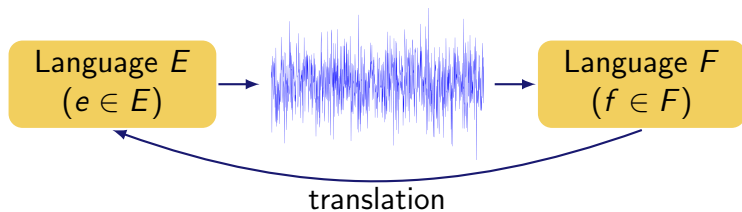
e : Good morning!

f : Bon jour!



SMT, basics

The Noisy Channel approach



Mathematically:

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

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

SMT, basics

Components

$$T(f) = \hat{e} = \operatorname{argmax}_e P(\textcolor{red}{e}) P(f|e)$$

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*

SMT, basics

Components

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

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*

SMT, basics

Components

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

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*

Outline

- 1 Introduction
- 2 Basics
- 3 Components**
 - Language model
 - Translation model
 - Decoder
- 4 The log-linear model
- 5 Beyond standard SMT
- 6 MT Evaluation

SMT, components

The language model $P(e)$

Language model

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

Estimation of how probable a sentence is.

Naïve estimation on a corpus with N sentences:

Frequentist probability
of a sentence e :

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

Problem:

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

SMT, components

The language model $P(e)$

Language model

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

Estimation of how probable a sentence is.

Naïve estimation on a corpus with N sentences:

Frequentist probability
of a sentence e :

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

Problem:

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

SMT, components

The language model $P(e)$

Language model

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

Estimation of how probable a sentence is.

Naïve estimation on a corpus with N sentences:

Frequentist probability
of a sentence e :

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

Problem:

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

SMT, components

The language model $P(e)$

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_i given the previous ones.
- Independence assumption: the probability of w_i is only conditioned by the n previous words.

SMT, components

The language model $P(e)$

Example, a 4-gram model

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

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

where, for each factor,

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

SMT, components

The language model $P(e)$

Example, a 4-gram model

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

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

where, for each factor,

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

SMT, components

The language model $P(e)$

Example, a 4-gram model

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

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

where, for each factor,

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

SMT, components

The language model $P(e)$

Example, a 4-gram model

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

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

where, for each factor,

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

SMT, components

The language model $P(e)$

Example, a 4-gram model

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

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

where, for each factor,

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

SMT, components

The language model $P(e)$

Example, a 4-gram model

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

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

where, for each factor,

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

SMT, components

The language model $P(e)$

Example, a 4-gram model

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

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

where, for each factor,

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

SMT, components

The language model $P(e)$

Example, a 4-gram model

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

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

where, for each factor,

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

SMT, components

The language model $P(e)$

Example, a 4-gram model

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

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

where, for each factor,

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

SMT, components

The language model $P(e)$

Example, a 4-gram model

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

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

where, for each factor,

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

SMT, components

The language model $P(e)$

Example, a 4-gram model

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

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

where, for each factor,

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

SMT, components

The language model $P(e)$

Example, a 4-gram model

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

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

where, for each factor,

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

SMT, components

The language model $P(e)$

Main problems and criticisms:

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

Solution

Smoothing techniques:

- Linear interpolation.

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

SMT, components

The language model $P(e)$

Main problems and criticisms:

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

Solution

Smoothing techniques:

- Linear interpolation.
- Back-off models.

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

SMT, components

The language model $P(e)$

Main problems and criticisms:

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

Solution

Smoothing techniques:

- Linear interpolation.

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

SMT, components

The language model $P(e)$

Main problems and criticisms:

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

Solution

Smoothing techniques:

- Linear interpolation.

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

SMT, components

The language model $P(e)$

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.

SMT, components

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?



SMT, components

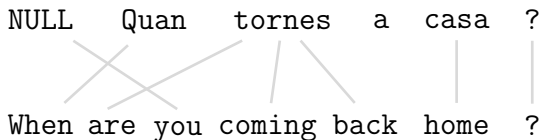
The translation model $P(f|e)$

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?



SMT, components

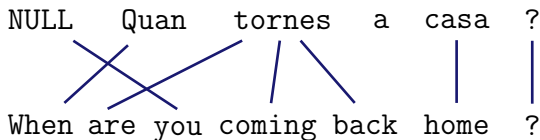
The translation model $P(f|e)$

Translation model

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

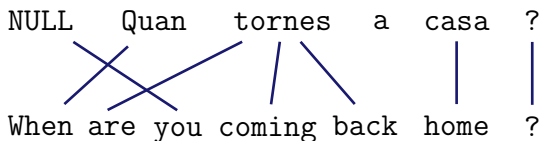
Estimation of the lexical correspondence between languages.

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



SMT, components

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.

SMT, components

The translation model $P(f|e)$

Word-based models: the IBM models

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

- Lexical probability t
 $t(\text{Quan}|\text{When})$: the prob. that **Quan** translates into **When**.
- Fertility n
 $n(3|\text{tornes})$: the prob. that **tornes** generates 3 words.

SMT, components

The translation model $P(f|e)$

Word-based models: the IBM models

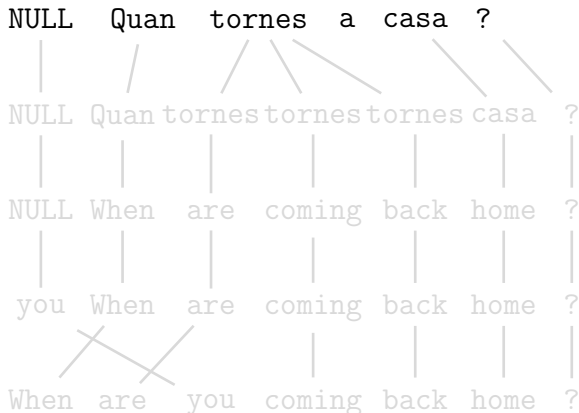
They characterise $P(f|e)$ with 4 parameters: t , n , d and p_1 .

- 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 p_1
 $p(\text{you}|\text{NULL})$: the prob. that the spurious word you is generated (from NULL).

SMT, components

The translation model $P(f|e)$

Back to the example:



Fertility

Translation

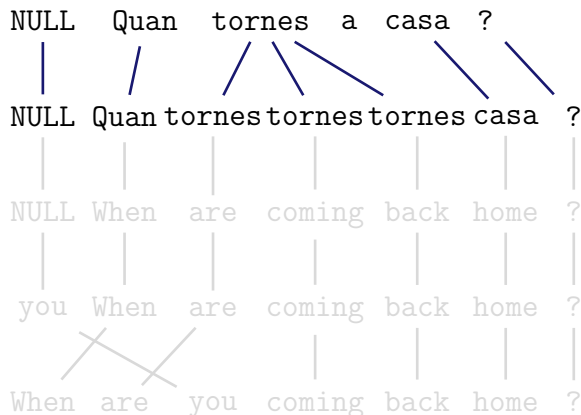
Insertion

Distortion

SMT, components

The translation model $P(f|e)$

Back to the example:



Fertility

Translation

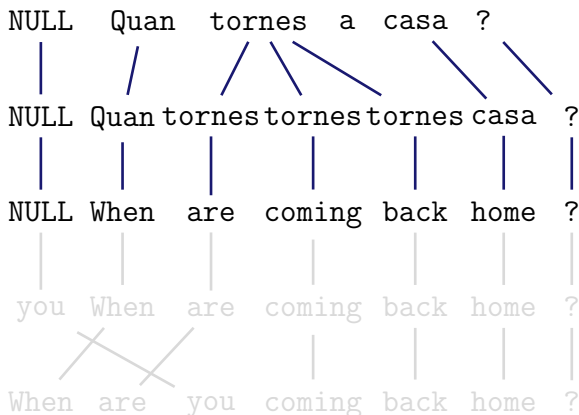
Insertion

Distortion

SMT, components

The translation model $P(f|e)$

Back to the example:



Fertility

Translation

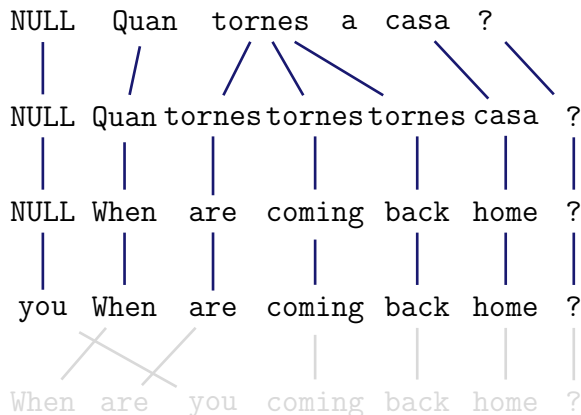
Insertion

Distortion

SMT, components

The translation model $P(f|e)$

Back to the example:



Fertility

Translation

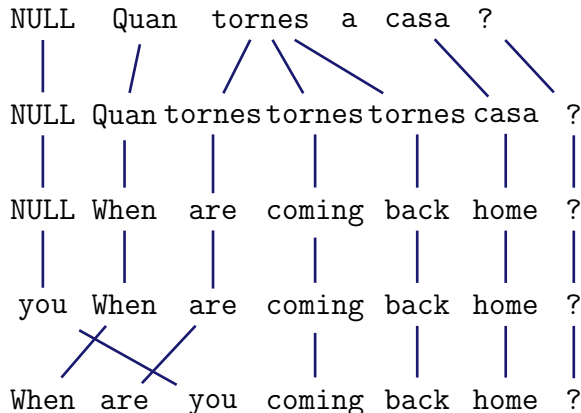
Insertion

Distortion

SMT, components

The translation model $P(f|e)$

Back to the example:



Fertility

Translation

Insertion

Distortion

SMT, components

The translation model $P(f|e)$

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.

SMT, components

The translation model $P(f|e)$

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.

SMT, components

The translation model $P(f|e)$

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.

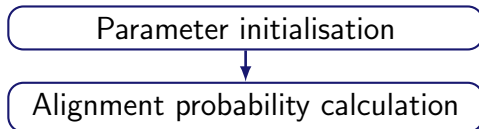
Solutions

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

SMT, components

The translation model $P(f|e)$

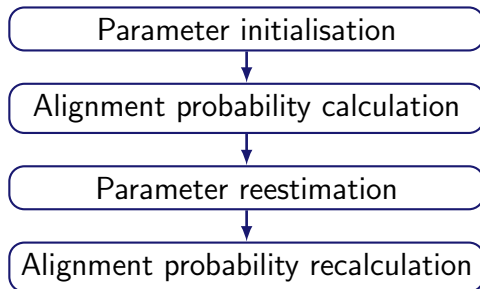
Expectation-Maximisation algorithm



SMT, components

The translation model $P(f|e)$

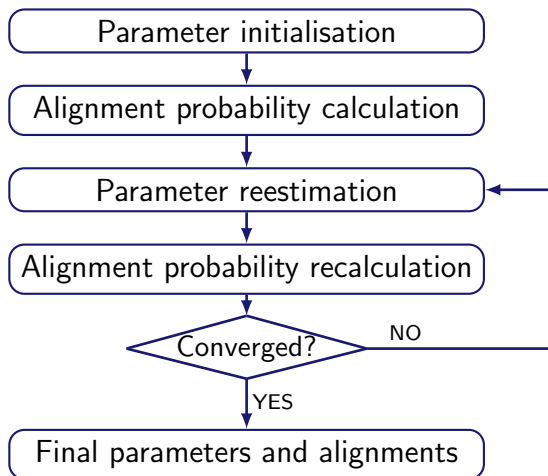
Expectation-Maximisation algorithm



SMT, components

The translation model $P(f|e)$

Expectation-Maximisation algorithm



SMT, components

The translation model $P(f|e)$

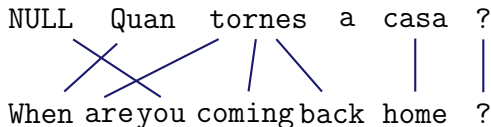
Alignment's asymmetry

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**.

Catalan
to
English

NULL Quan tornes a casa ?
When are you coming back home ?



English
to
Catalan

NULL When are you coming back home ?
Quan tornes a casa ?



SMT, components

The translation model $P(f|e)$

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**.

Catalan
to
English

NULL Quan tornes a casa ?
When are you coming back home ?

English
to
Catalan

NULL When are you coming back home ?
Quan tornes a casa ?

SMT, components

The translation model $P(f|e)$

Graphically:

	NULL	Quan	tornes	a	casa	?
NULL						
When						
are						
you						
coming						
back						
home						
?						

Catalan to English

SMT, components

The translation model $P(f|e)$

Graphically:

	NULL	Quan	tornes	a	casa	?
NULL						
When						
are						
you						
coming						
back						
home						
?						

English to Catalan

SMT, components

The translation model $P(f|e)$

Alignment symmetrisation

- Intersection: high-confidence, high precision.

	NULL	Quan	tornes	a	casa	?
NULL						
When						
are						
you						
coming						
back						
home						
?						

Catalan to English \cap English to Catalan

SMT, components

The translation model $P(f|e)$

Alignment symmetrisation

- Union: lower confidence, high recall.

	NULL	Quan	tornes	a	casa	?
NULL						
When						
are						
you						
coming						
back						
home						
?						

Catalan to English \cup English to Catalan

SMT, components

The translation model $P(f|e)$

From Word-based to Phrase-based models

f: En David llegeix el llibre nou.

SMT, components

The translation model $P(f|e)$

From Word-based to Phrase-based models

f: En David llegeix el llibre nou.

e: ϕ

SMT, components

The translation model $P(f|e)$

From Word-based to Phrase-based models

f: En **David** llegeix el llibre nou.

e: **David**

SMT, components

The translation model $P(f|e)$

From Word-based to Phrase-based models

f: En David **llegeix** el llibre nou.

e: David **reads**

SMT, components

The translation model $P(f|e)$

From Word-based to Phrase-based models

f: En David llegeix el llibre nou.

e: David reads the

SMT, components

The translation model $P(f|e)$

From Word-based to Phrase-based models

f: En David llegeix el **llibre** nou.

e: David reads the **book**

SMT, components

The translation model $P(f|e)$

From Word-based to Phrase-based models

f: En David llegeix el llibre nou.

e: David reads the book new.

SMT, components

The translation model $P(f|e)$

From Word-based to Phrase-based models

f: En David llegeix el llibre nou.

e: David reads the book new. ~

SMT, components

The translation model $P(f|e)$

From Word-based to Phrase-based models

f: En David llegeix el llibre nou.

e: David reads the new book. ✓

SMT, components

The translation model $P(f|e)$

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.

SMT, components

The translation model $P(f|e)$

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: ϕ

SMT, components

The translation model $P(f|e)$

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**

SMT, components

The translation model $P(f|e)$

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

SMT, components

The translation model $P(f|e)$

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

SMT, components

The translation model $P(f|e)$

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**

SMT, components

The translation model $P(f|e)$

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

SMT, components

The translation model $P(f|e)$

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.

SMT, components

The translation model $P(f|e)$

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. ✗

SMT, components

The translation model $P(f|e)$

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: ϕ

SMT, components

The translation model $P(f|e)$

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

SMT, components

The translation model $P(f|e)$

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

SMT, components

The translation model $P(f|e)$

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

SMT, components

The translation model $P(f|e)$

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

SMT, components

The translation model $P(f|e)$

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.

SMT, components

The translation model $P(f|e)$

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. ✓

SMT, components

The translation model $P(f|e)$

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.

SMT, components

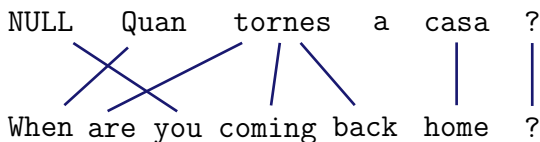
The translation model $P(f|e)$

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 caught.

SMT, components

The translation model $P(f|e)$

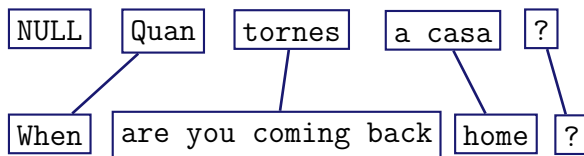


With the new translation units, $P(f|e)$ can be obtained following the **same strategy** as for word-based models with few modifications:

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

SMT, components

The translation model $P(f|e)$

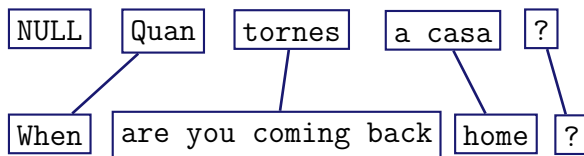


With the new translation units, $P(f|e)$ can be obtained following the **same strategy** as for word-based models with few modifications:

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

SMT, components

The translation model $P(f|e)$

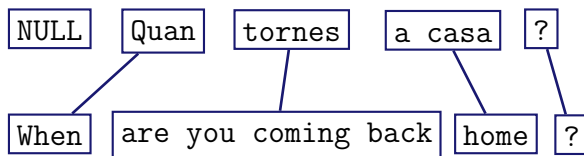


With the new translation units, $P(f|e)$ can be obtained following the **same strategy** as for word-based models with few modifications:

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

SMT, components

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?

SMT, components

The translation model $P(f|e)$

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.

We do not use the term phrase here in its linguistic sense: a phrase can be any sequence of words, even if they are not a linguistic constituent.

SMT, components

The translation model $P(f|e)$

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.

We do not use the term phrase here in its linguistic sense: a phrase can be any sequence of words, even if they are not a linguistic constituent.

SMT, components

The translation model $P(f|e)$

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.

We do not use the term phrase here in its linguistic sense: a phrase can be any sequence of words, even if they are not a linguistic constituent.

SMT, components

The translation model $P(f|e)$

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.¹

We do not use the term phrase here in its linguistic sense: a phrase can be any sequence of words, even if they are not a linguistic constituent.

SMT, components

The translation model $P(f|e)$

Phrase extraction through an example:

	Quan	tornes	tu	a	casa	?
When						
are						
you						
coming						
back						
home						
?						

(Quan tornes, When are you coming back)

SMT, components

The translation model $P(f|e)$

Phrase extraction through an example:

	Quan	tornes	tu	a	casa	?
When						
are						
you						
coming						
back						
home						
?						

~~(Quan tornes, When are you coming back)~~

SMT, components

The translation model $P(f|e)$

Phrase extraction through an example:

	Quan	tornes	tu	a	casa	?
When						
are						
you						
coming						
back						
home						
?						

~~(Quan tornes, When are you coming back)~~

(Quan tornes tu, When are you coming back)

SMT, components

The translation model $P(f|e)$

Intersection

	Quan	tornes	a	casa	?
When					
are					
you					
coming					
back					
home					
?					

(Quan, When) (Quan tornes, When are you coming) (Quan tornes a casa, When are you coming back home) (Quan tornes a casa ?, When are you coming back home ?) (tornes, coming) (tornes a casa, coming back home) (tornes a casa ?, coming back home ?) (casa, home) (casa ?, home ?) (?, ?) 10 phrases

SMT, components

The translation model $P(f|e)$

Intersection

	Quan	tornes	a	casa	?
When					
are					
you					
coming					
back					
home					
?					

(Quan, When) (Quan tornes, When are you coming) (Quan tornes a casa, When are you coming back home) (Quan tornes a casa ?, When are you coming back home ?) (tornes, coming) (tornes a casa, coming back home) (tornes a casa ?, coming back home ?) (casa, home) (casa ?, home ?) (?, ?) 10 phrases

SMT, components

The translation model $P(f|e)$

Intersection

	Quan	tornes	a	casa	?
When					
are					
you					
coming					
back					
home					
?					

(Quan, When) (Quan tornes, When are you coming) (Quan tornes a casa, When are you coming back home) (Quan tornes a casa ?, When are you coming back home ?) (tornes, coming) (tornes a casa, coming back home) (tornes a casa ?, coming back home ?) (casa, home) (casa ?, home ?) (?, ?) 10 phrases

SMT, components

The translation model $P(f|e)$

Intersection

	Quan	tornes	a	casa	?
When					
are					
you					
coming					
back					
home					
?					

(Quan, When) (Quan tornes, When are you coming) (Quan tornes a casa, When are you coming back home) (Quan tornes a casa ?, When are you coming back home ?) (tornes, coming) (tornes a casa, coming back home) (tornes a casa ?, coming back home ?) (casa, home) (casa ?, home ?) (?, ?) 10 phrases

SMT, components

The translation model $P(f|e)$

Intersection

	Quan	tornes	a	casa	?
When					
are					
you					
coming					
back					
home					
?					

(Quan, When) (Quan tornes, When are you coming) (Quan tornes a casa, When are you coming back home) (Quan tornes a casa ?, When are you coming back home ?) (tornes, coming) (tornes a casa, coming back home) (tornes a casa ?, coming back home ?) (casa, home) (casa ?, home ?) (?, ?) 10 phrases

SMT, components

The translation model $P(f|e)$

Intersection

	Quan	tornes	a	casa	?
When					
are					
you					
coming					
back					
home					
?					

(Quan, When) (Quan tornes, When are you coming) (Quan tornes a casa, When are you coming back home) (Quan tornes a casa ?, When are you coming back home ?) (tornes, coming) (tornes a casa, coming back home) (tornes a casa ?, coming back home ?) (casa, home) (casa ?, home ?) (?, ?) 10 phrases

SMT, components

The translation model $P(f|e)$

Intersection

	Quan	tornes	a	casa	?
When					
are					
you					
coming					
back					
home					
?					

(Quan, When) (Quan tornes, When are you coming) (Quan tornes a casa, When are you coming back home) (Quan tornes a casa ?, When are you coming back home ?) (tornes, coming) (tornes a casa, coming back home) (tornes a casa ?, coming back home ?) (casa, home) (casa ?, home ?) (?, ?) 10 phrases

SMT, components

The translation model $P(f|e)$

Intersection

	Quan	tornes	a	casa	?
When					
are					
you					
coming					
back					
home					
?					

(Quan, When) (Quan tornes, When are you coming) (Quan tornes a casa, When are you coming back home) (Quan tornes a casa ?, When are you coming back home ?) (tornes, coming) (tornes a casa, coming back home) (tornes a casa ?, coming back home ?) (casa, home) (casa ?, home ?) (?, ?) 10 phrases

SMT, components

The translation model $P(f|e)$

Intersection

	Quan	tornes	a	casa	?
When					
are					
you					
coming					
back					
home					
?					

(Quan, When) (Quan tornes, When are you coming) (Quan tornes a casa, When are you coming back home) (Quan tornes a casa ?, When are you coming back home ?) (tornes, coming) (tornes a casa, coming back home) (tornes a casa ?, coming back home ?) (casa, home) (casa ?, home ?) (?, ?) 10 phrases

SMT, components

The translation model $P(f|e)$

Intersection

	Quan	tornes	a	casa	?
When					
are					
you					
coming					
back					
home					
?					

(Quan, When) (Quan tornes, When are you coming) (Quan tornes a casa, When are you coming back home) (Quan tornes a casa ?, When are you coming back home ?) (tornes, coming) (tornes a casa, coming back home) (tornes a casa ?, coming back home ?) (casa, home) (casa ?, home ?) (?, ?) 10 phrases

SMT, components

The translation model $P(f|e)$

Union

	Quan	tornes	a	casa	?
When					
are					
you					
coming					
back					
home					
?					

(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

SMT, components

The translation model $P(f|e)$

Union

	Quan	tornes	a	casa	?
When					
are					
you					
coming					
back					
home					
?					

(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

SMT, components

The translation model $P(f|e)$

Union

	Quan	tornes	a	casa	?
When					
are					
you					
coming					
back					
home					
?					

(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

SMT, components

The translation model $P(f|e)$

Union

	Quan	tornes	a	casa	?
When					
are					
you					
coming					
back					
home					
?					

(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

SMT, components

The translation model $P(f|e)$

Union

	Quan	tornes	a	casa	?
When					
are					
you					
coming					
back					
home					
?					

(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**

SMT, components

The translation model $P(f|e)$

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.

SMT, components

The translation model $P(f|e)$

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**.

SMT, components

The translation model $P(f|e)$

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.

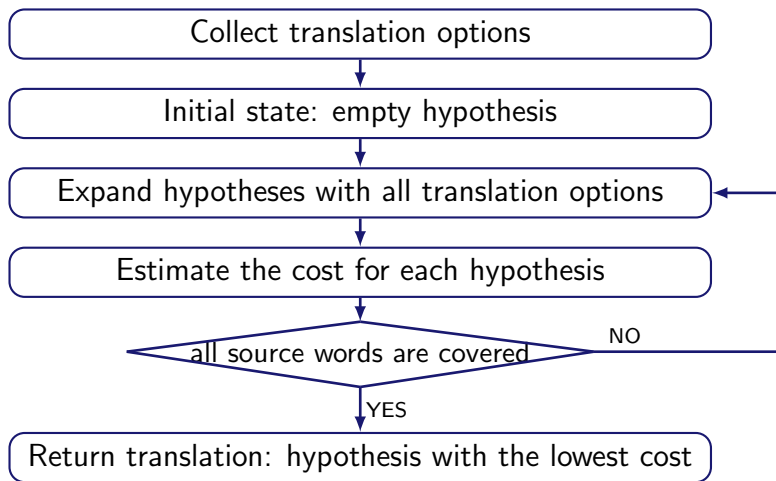
In our context, one can find:

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

SMT, components

A beam-search decoder

Core algorithm



SMT, components

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)

SMT, components

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 - -

SMT, components

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** - -

SMT, components

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:

ϕ

Source words already translated:

- - - -

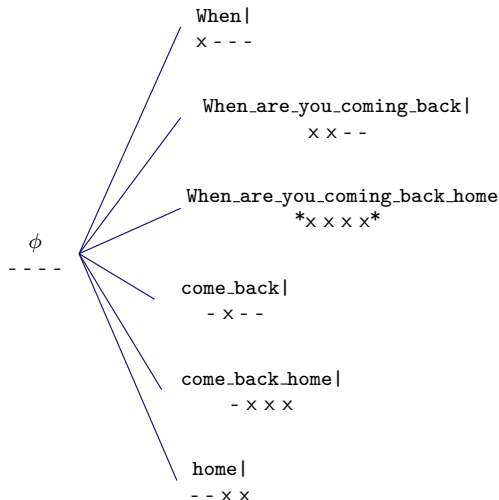
SMT, components

A beam-search decoder

ϕ

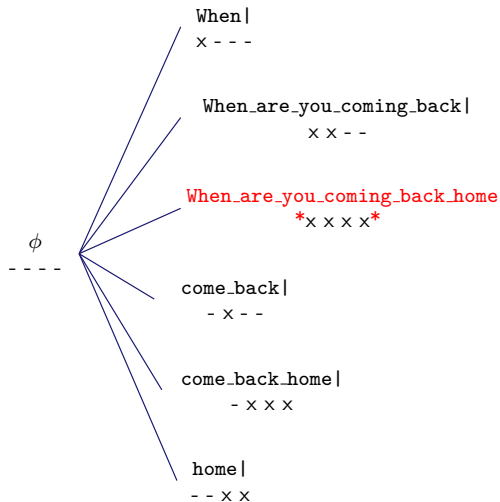
SMT, components

A beam-search decoder



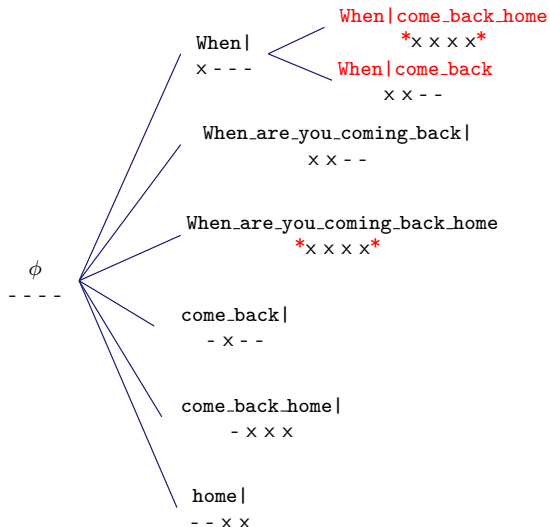
SMT, components

A beam-search decoder



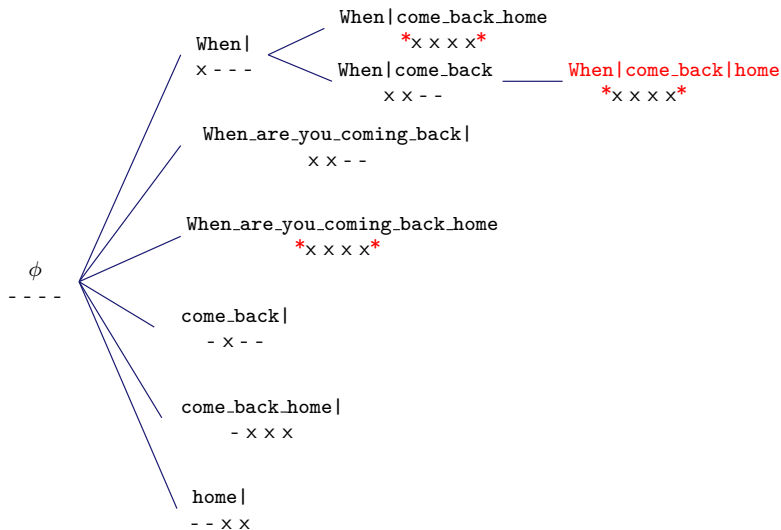
SMT, components

A beam-search decoder



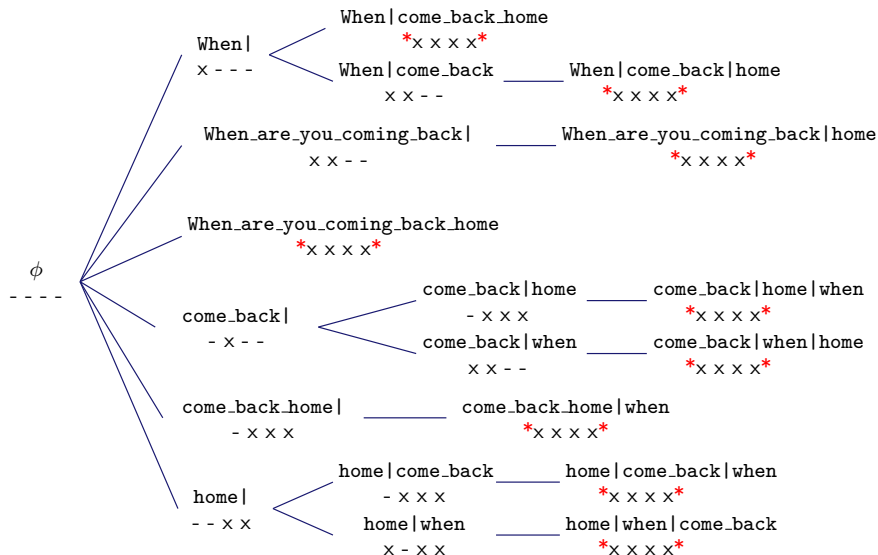
SMT, components

A beam-search decoder



SMT, components

A beam-search decoder



SMT, components

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

SMT, components

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

SMT, components

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

SMT, components

A beam-search decoder

Hypotheses recombination

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

When | come_back_home \longleftrightarrow When | come_back | home
x x x x x x x x

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

SMT, components

A beam-search decoder

Hypotheses recombination

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

When | come_back_home \iff When | come_back | home
x x x x x x x x

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

SMT, components

A beam-search decoder

Hypotheses recombination

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

When | come_back_home \iff When | come_back | home
x x x x x x x x

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

SMT, components

A beam-search decoder

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**.

SMT, components

A beam-search decoder

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.

SMT, components

Decoder

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.

Outline

- 1 Introduction
- 2 Basics
- 3 Components
- 4 The log-linear model**
- 5 Beyond standard SMT
- 6 MT Evaluation

SMT, the log-linear model

Motivation

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

SMT, the log-linear model

Motivation

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

SMT, the log-linear model

Motivation

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

SMT, the log-linear model

Motivation

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

$$h_1(f, e) = \log P(e), \quad h_2(f, e) = \log P(f|e), \quad \text{and} \quad \lambda_1 = \lambda_2 = 1$$

\Rightarrow Maximum likelihood model

SMT, the log-linear model

Motivation

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** λ_m must be somehow determined.
- Different knowledge sources can be used.

SMT, the log-linear model

Features

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.

SMT, the log-linear model

Features

State of the art feature functions

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

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

SMT, the log-linear model

Features

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 $P_d(e, f)$
 $P_d(\text{ini}_{\text{phrase}_i}, \text{end}_{\text{phrase}_{i-1}})$: Relative distortion probability distribution. A simple distortion model:
$$P_d(\text{ini}_{\text{phrase}_i}, \text{end}_{\text{phrase}_{i-1}}) = \alpha |\text{ini}_{\text{phrase}_i} - \text{end}_{\text{phrase}_{i-1}} - 1|$$

SMT, the log-linear model

Weights optimisation

Development training, weights optimisation

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

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).

SMT, the log-linear model

Weights optimisation

Development training, weights optimisation

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

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).

SMT, the log-linear model

Weights optimisation

Development training, weights optimisation

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

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).**

SMT, the log-linear model

Minimum Error-Rate Training (MERT)

Minimum Error-Rate Training

- Approach: Minimise an error function.

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

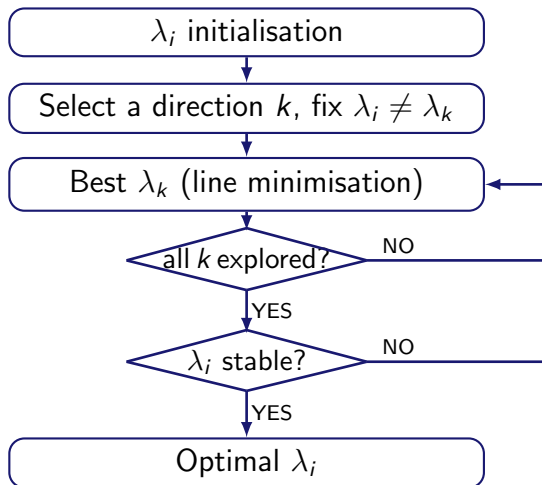
- There exist several error measures or metrics.
- Metrics not always correlate with human judgements.
- The quality of the final translation on the metric chosen for the optimisation is shown to improve.
- For the moment, let's say we use BLEU.

(More on MT Evaluation section)

SMT, the log-linear model

Minimum Error-Rate Training (MERT)

Minimum Error-Rate Training rough algorithm



SMT, the log-linear model

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.

Outline

- 1 Introduction
- 2 Basics
- 3 Components
- 4 The log-linear model
- 5 Beyond standard SMT**
 - Factored translation models
 - Syntactic translation models
 - Ongoing research
- 6 MT Evaluation

SMT, beyond standard SMT

Including linguistic information

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.

SMT, beyond standard SMT

Factored translation models

Factored translation models

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

$$(\text{word}) \implies (\text{word}, \text{lemma}, \text{PoS}, \text{morphology}, \dots)$$

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

SMT, beyond standard SMT

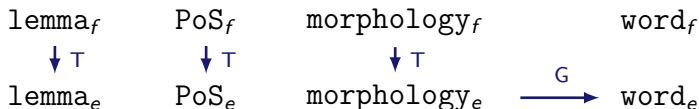
Factored translation models

Factored translation models

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

$$(\text{word}) \implies (\text{word}, \text{lemma}, \text{PoS}, \text{morphology}, \dots)$$

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



SMT, beyond standard SMT

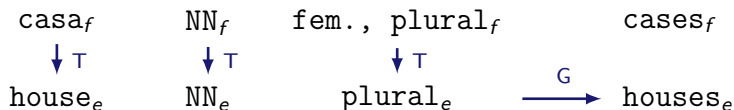
Factored translation models

Factored translation models

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

$$(\text{word}) \implies (\text{word}, \text{lemma}, \text{PoS}, \text{morphology}, \dots)$$

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



SMT, beyond standard SMT

Factored translation models

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.

SMT, beyond standard SMT

Syntactic translation models

Syntactic translation models

Incorporate syntax to the source and/or target languages.

Approaches

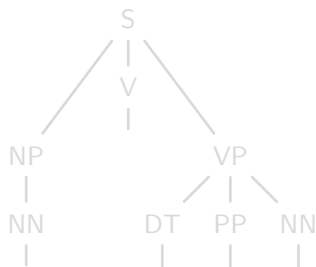
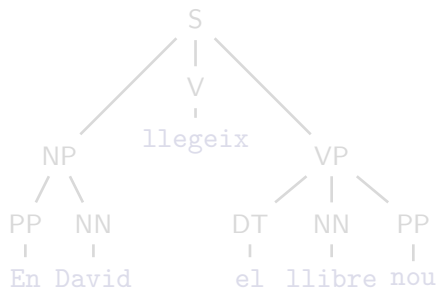
- Synchronous grammar formalism which learns a grammar that can simultaneously generate both trees.
 - ▶ **Syntax-based**. Respect linguistic units in translation.
 - ▶ **Hierarchical phrase-based**. Respect phrases in translation.

SMT, beyond standard SMT

Syntax-based translation models

Syntactic models ease reordering. An intuitive example:

En David llegeix un llibre nou

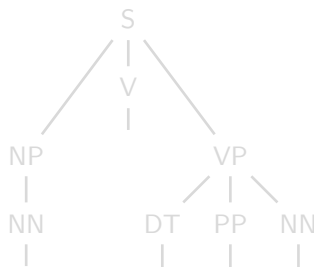
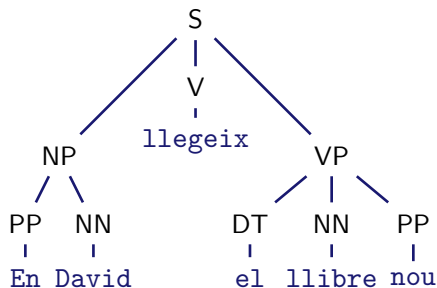


SMT, beyond standard SMT

Syntax-based translation models

Syntactic models ease reordering. An intuitive example:

En David llegeix un llibre nou

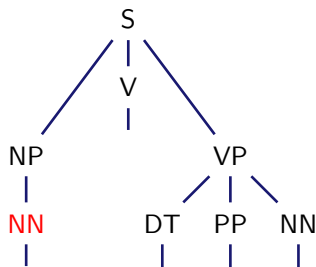
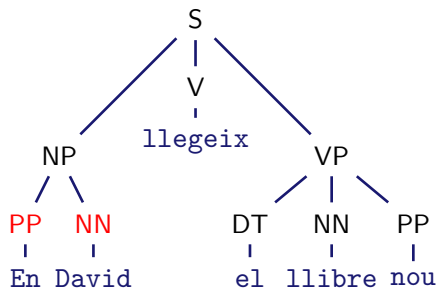


SMT, beyond standard SMT

Syntax-based translation models

Syntactic models ease reordering. An intuitive example:

En David llegeix un llibre nou

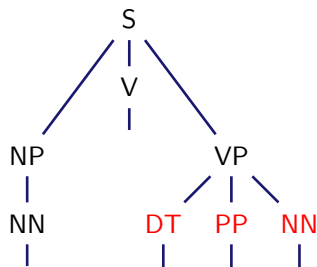
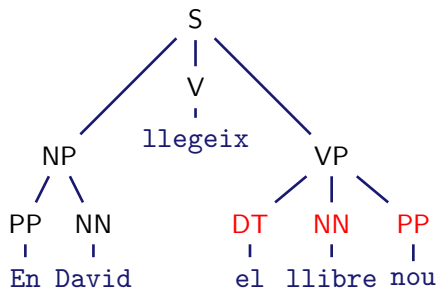


SMT, beyond standard SMT

Syntax-based translation models

Syntactic models ease reordering. An intuitive example:

En David llegeix un llibre nou

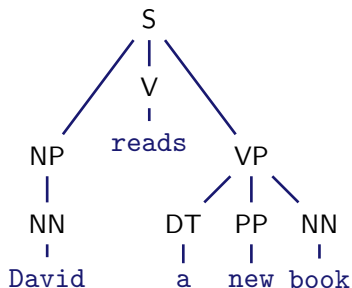
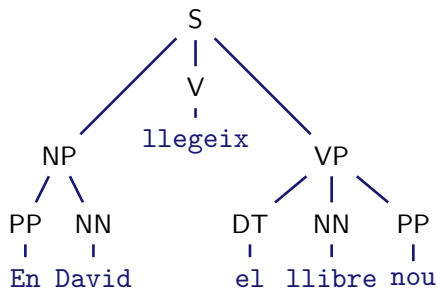


SMT, beyond standard SMT

Syntax-based translation models

Syntactic models ease reordering. An intuitive example:

En David llegeix un llibre nou



David reads a new book

SMT, beyond standard SMT

Ongoing research

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.

SMT, beyond standard SMT

Ongoing research

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.

SMT, beyond standard SMT

Ongoing research

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

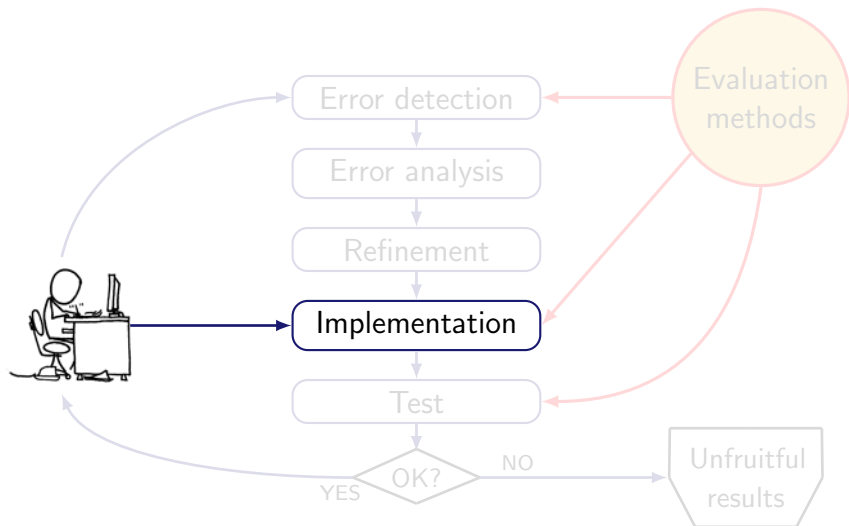
- Factored models include linguistic information in phrase-based models and are suitable for morphologically rich languages.
- Syntactic models consider somehow syntax and are adequate for language pairs with a different structure of the sentences.
- Current research addresses both new models and modifications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

Outline

- 1 Introduction
- 2 Basics
- 3 Components
- 4 The log-linear model
- 5 Beyond standard SMT
- 6 MT Evaluation**

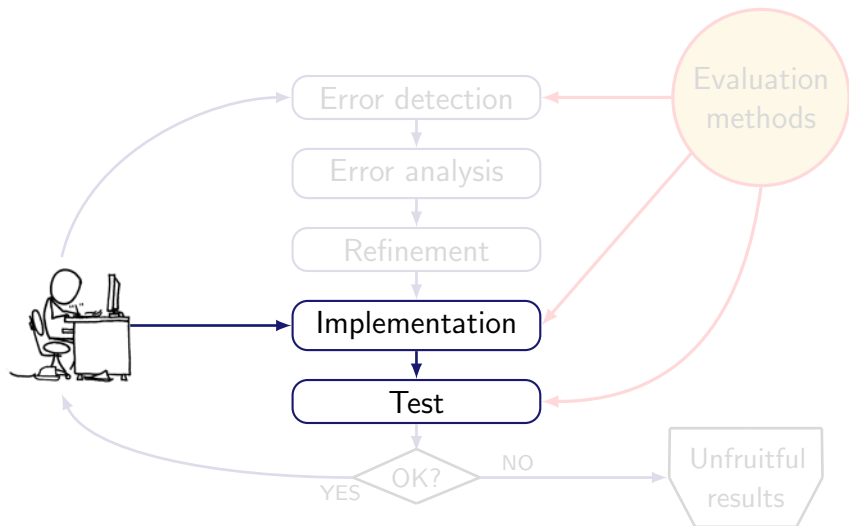
MT Evaluation

Importance for system development



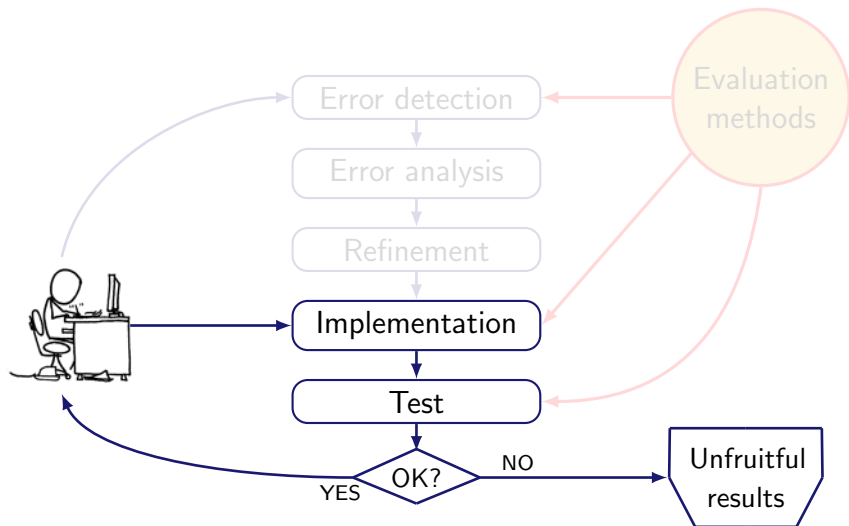
MT Evaluation

Importance for system development



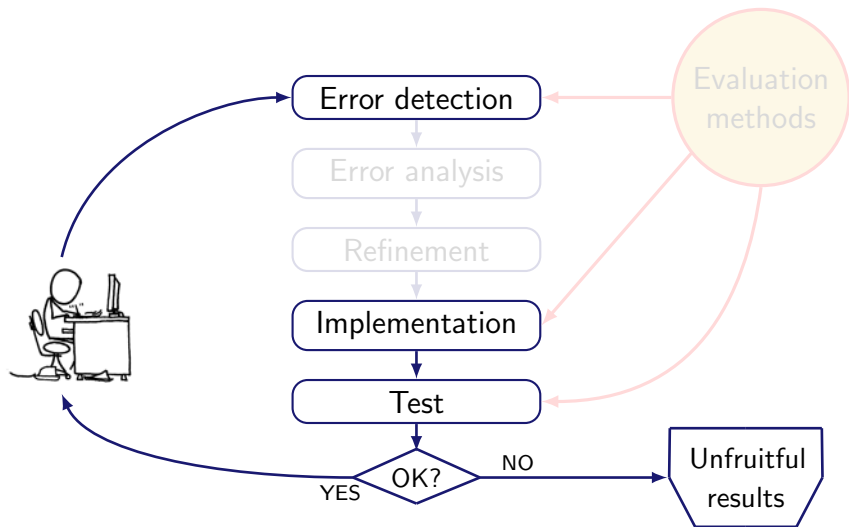
MT Evaluation

Importance for system development



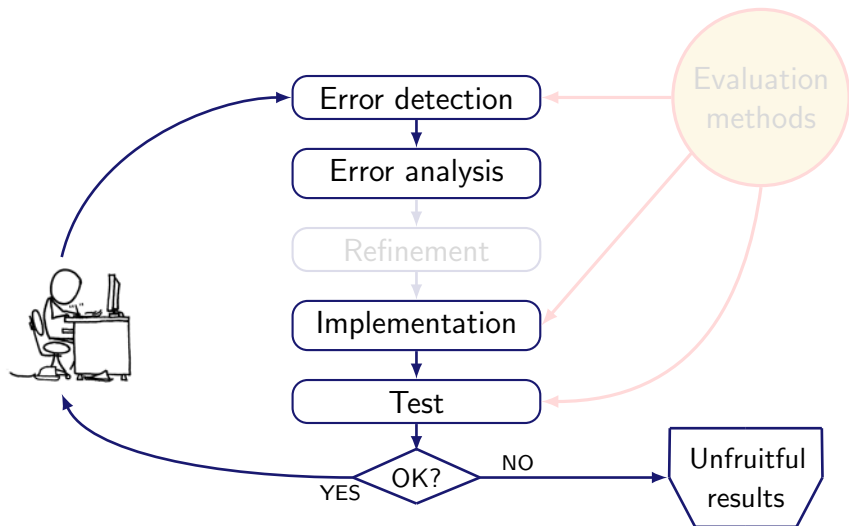
MT Evaluation

Importance for system development



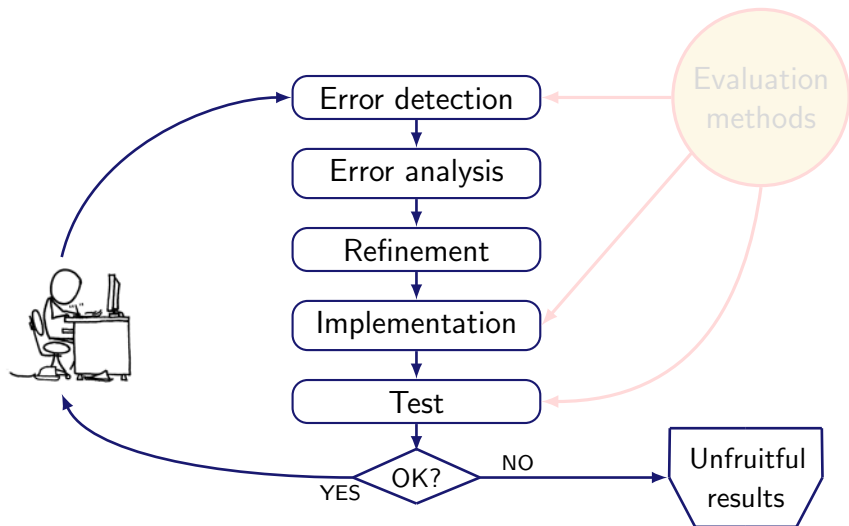
MT Evaluation

Importance for system development



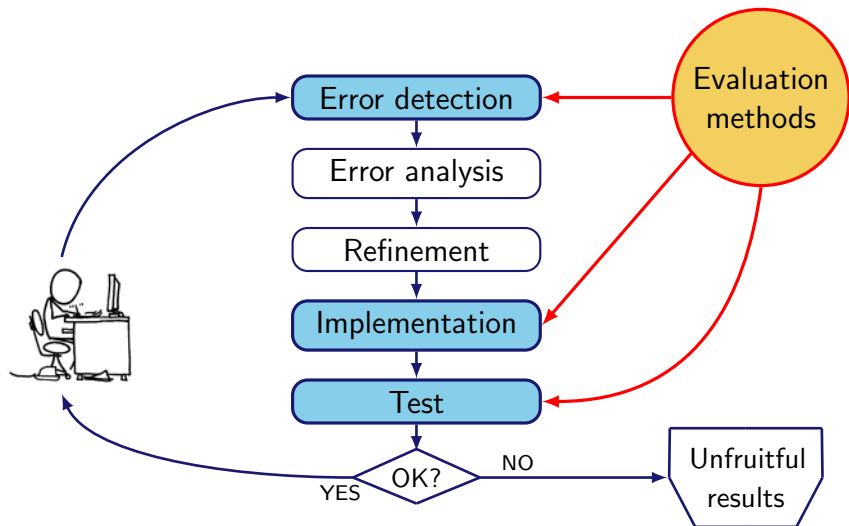
MT Evaluation

Importance for system development



MT Evaluation

Importance for system development



MT Evaluation

Automatic evaluation

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)

MT Evaluation

Lexical similarity

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

MT Evaluation

Lexical similarity

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.

MT Evaluation

Lexical similarity

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 translation **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 neither a sufficient nor a necessary condition so that two sentences convey the same meaning.

MT Evaluation

Lexical similarity

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 translation 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 neither a sufficient nor a necessary condition so that two sentences convey the same meaning.

MT Evaluation

Lexical similarity

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 translation 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.

MT Evaluation

Ongoing research

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.

MT Evaluation

Ongoing research

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.

MT Evaluation

Ongoing research

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.

MT Evaluation

Summary

MT Evaluation: keep in mind

- Evaluation is important in the system development cycle. Automatic evaluation accelerates significantly 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.

Thanks!

A last alignment

Gràcies a

en Jesús Giménez

per

algunes transparències

Thanks to

Jesús Giménez

for

some of the material



?

Part II

SMT experiments

7 Translation system

- Software
- Steps

8 Evaluation system

- Software
- Steps

Build your own SMT system

- 1 Language model with SRILM.
<http://www.speech.sri.com/projects/srilm/download.htm>
- 2 Word alignments with GIZA++.
<http://code.google.com/p/giza-pp/downloads/list>
- 3 And everything else with the Moses package.
<http://sourceforge.net/projects/mosesdecoder>

1. Download and prepare your data

- ➊ Parallel corpora and some tools can be downloaded for instance from the WMT 2010 web page:
<http://www.statmt.org/wmt10/translation-task.html>

How to construct a baseline system is also explained there:
<http://www.statmt.org/wmt10/baseline.html>

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

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

- 2 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)

- 3 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
```

- 4 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

- 1 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
```

3. Train the translation model

- 1 Use the Moses script `train-factored-phrase-model.perl`

This script performs the whole training:

```
cristina@cosmos:~$ train-factored-phrase-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)

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

```
bin/moses-scripts/training/train-factored-phrase-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

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

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

- 2 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

- 1 Tokenise and lowercase the test set as before.
- 2 Filter the model with Moses script.
(mandatory for large translation tables)

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

- 3 Run the decoder:

```
moses/moses-cmd/src/moses -config ./filteredmodel/moses.ini  
-input-file testset.es > testset.translated.en
```

Evaluate the results

- 1 With BLEU scoring tool. Available as a Moses script or from NIST:

<http://www.itl.nist.gov/iad/mig/tools/mtevalv13a-20091001.tar.gz>

- 2 With IQmt package.

<http://www.lsi.upc.edu/~nlp/IQMT/>

MT Evaluation

Steps

1. Evaluate the results

- 1 With BLEU scoring tool in Moses:

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

MT Evaluation

Steps

2. Evaluate the results on-line

- ① OpenMT Evaluation Demo

<http://biniki.lsi.upc.edu/openMT/evaldemo.php>

Part III

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^a 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 International 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.



Jessí 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.



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.

References VIII



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.

References IX



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.