

# A proposal for an Arabic-to-English SMT system

Cristina España i Bonet

Advisor: Dr. Lluís Màrquez Villodre

22th February, 2008

- 1 Introduction
  - Statistical Machine Translation
  - Language Pair
- 2 System Design
- 3 Experiments and evaluation
- 4 Conclusions

# Introduction

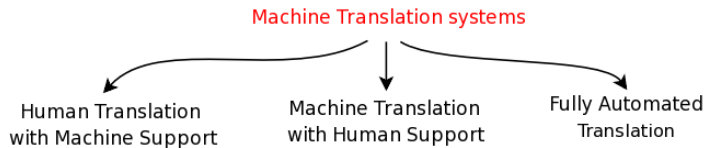
## Machine Translation

**Goal:** *Apply Machine Translation techniques to translate from Arabic to English in the context of the 2008 NIST Machine Translation Open evaluation*

# Introduction

## Machine Translation

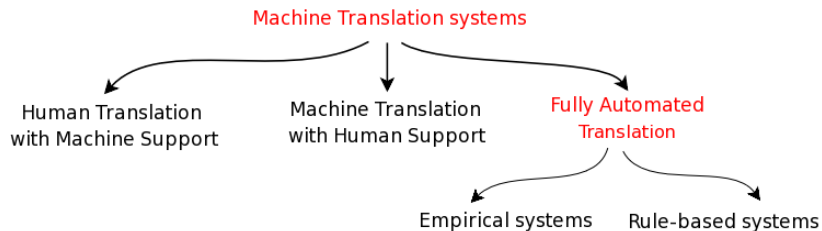
**Goal:** *Apply Machine Translation techniques to translate from Arabic to English in the context of the 2008 NIST Machine Translation Open evaluation*



# Introduction

## Machine Translation

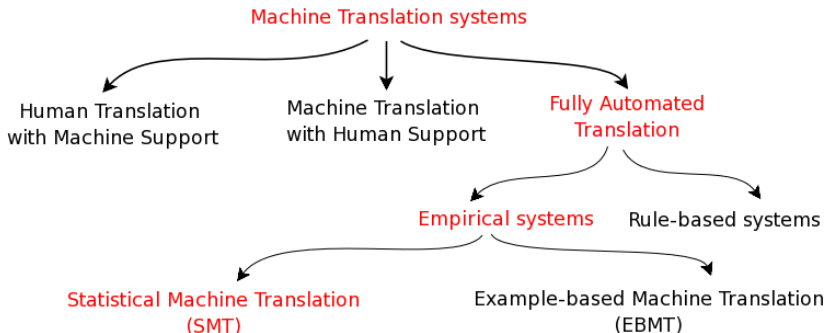
**Goal:** *Apply Machine Translation techniques to translate from Arabic to English in the context of the 2008 NIST Machine Translation Open evaluation*



# Introduction

## Machine Translation

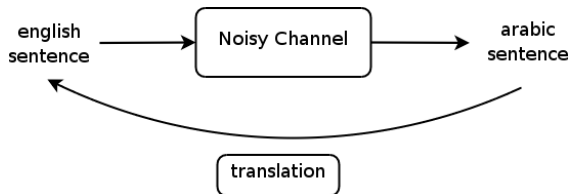
**Goal:** *Apply Machine Translation techniques to translate from Arabic to English in the context of the 2008 NIST Machine Translation Open evaluation*



# Statistical Machine Translation

## Translation as a Noisy Channel Model

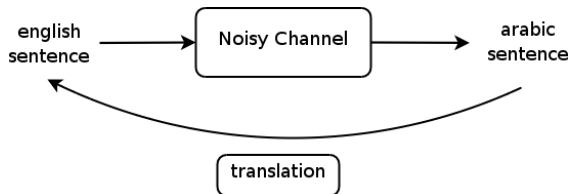
Basics:



# Statistical Machine Translation

## Translation as a Noisy Channel Model

Basics:



Mathematically:

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

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



# Statistical Machine Translation

## Components

### $P(e)$ Language Model

- Probability scores of n-grams in the target language
- Takes care of correctness and fluency
- Data: corpora in the target language

# Statistical Machine Translation

## Components

### $P(e)$ Language Model

- Probability scores of n-grams in the target language
- Takes care of correctness and fluency
- Data: corpora in the target language

### $P(f|e)$ Translation Model

- IBM models
- Data: aligned corpora in source and target languages

# Statistical Machine Translation

## Components

### $P(e)$ Language Model

- Probability scores of n-grams in the target language
- Takes care of correctness and fluency
- Data: corpora in the target language

### $P(f|e)$ Translation Model

- IBM models
- Data: aligned corpora in source and target languages

### argmax

- Search done by the *decoder*

# Statistical Machine Translation

Log-linear model

Maximum Likelihood estimate  $\leftrightarrow$  Maximum Entropy estimate

# Statistical Machine Translation

## Log-linear model

Maximum Likelihood estimate  $\leftrightarrow$  Maximum Entropy estimate



$$\operatorname{argmax} \log P(e|f) = \operatorname{argmax} \sum_m \lambda_m h_m(f|e)$$

# Statistical Machine Translation

## Log-linear model

Maximum Likelihood estimate  $\leftrightarrow$  Maximum Entropy estimate



$$\operatorname{argmax} \log P(e|f) = \operatorname{argmax} \sum_m \lambda_m h_m(f|e)$$

- 👉  $h_m \rightarrow$  features (log-probabilities)
  - 👉 Language and translation models,
  - 👉 and distortion, word penalty, phrase penalty...
- $\lambda_m \rightarrow$  weight of every feature

# Statistical Machine Translation

## Log-linear model

Maximum Likelihood estimate  $\leftrightarrow$  Maximum Entropy estimate



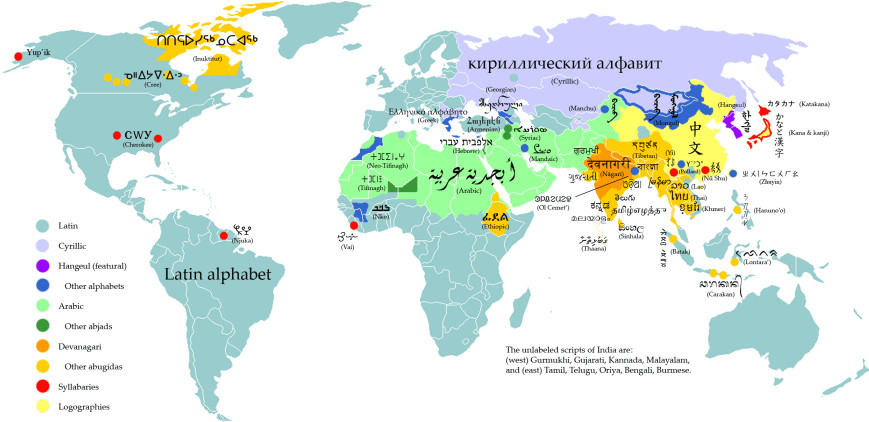
$$\operatorname{argmax} \log P(e|f) = \operatorname{argmax} \sum_m \lambda_m h_m(f|e)$$

- $h_m \rightarrow$  features (log-probabilities)
  - ▶ Language and translation models,
  - ▶ and distortion, word penalty, phrase penalty...

☞  $\lambda_m \rightarrow$  weight of every feature

# Language pair

## Arabic-English





# Language pair

Arabic

وتأتي دول فرنسا وبريطانيا وإيطاليا وألمانيا وأيرلندا وأسبانيا ولو كسمبورج في المقدمة وبخلاف الشركات الأوروبية فقد وصل حجم رؤوس الأموال المصدرة لل شركات العاملة في مصر حتي ديسمبر ٢٠٠٠ الي ١٢٦ مليار ج نيه لعدد ١٠ آلاف شركة استثمارية مؤسسه وفقا لقانون الاستثمار ..

☞ Right to left text (numerals: left to right)

- VSO structure
- Alphabet: allographic variants, diacritics and ligatures
- Agglutinative language

# Language pair

Arabic

وتأتي دول فرنسا وبريطانيا وإيطاليا وألمانيا وأيرلندا وأسبانيا ولو كسمبورج في المقدمة وبخلاف الشركات الأوروبية فقد وصل حجم رءوس الأموال المصدرة لل شركات العاملة في مصر حتي ديسمبر ٢٠٠٠ الي ١٢٦ مليار ج نيه لعدد ١٠ آلاف شركة استثمارية مؤسسه وفقا لقانون الاستثمار ..

- Right to left text (numerals: left to right)

🗨️ VSO structure

- Alphabet: allographic variants, diacritics and ligatures
- Agglutinative language

# Language pair

Arabic

وتأتي دول فرنسا وبريطانيا وإيطاليا وألمانيا وأيرلندا وأسبانيا ولو كسمبورج في المقدمة وبخلاف الشركات الأوروبية فقد وصل حجم رؤوس الأموال المصدرة لل شركات العاملة في مصر حتي ديسمبر ٢٠٠٠ الي ١٢٦ مليار ج نيه لعدد ١٠ آلاف شركة استثمارية مؤسسه وفقا لقانون الاستثمار ..

- Right to left text (numerals: left to right)
- VSO structure
- ☞ Alphabet: allographic variants, diacritics and ligatures
- Agglutinative language

# Language pair

Arabic

وتأتي دول فرنسا وبريطانيا وإيطاليا وألمانيا وأيرلندا وأسبانيا ولو كسمبورج في المقدمة وبخلاف الشركات الأوروبية فقد وصل حجم رؤوس الأموال المصدرة لل شركات العاملة في مصر حتي ديسمبر ٢٠٠٠ الي ١٢٦ مليار ج نيه لعدد ١٠ آلاف شركة استثمارية مؤسسه وفقا لقانون الاستثمار ..

- Right to left text (numerals: left to right)
- VSO structure
- Alphabet: allographic variants, diacritics and ligatures
- 👉 Agglutinative language

# Language pair

Arabic: diacritics

|           |     |     |       |       |    |     |       |       |       |     |    |
|-----------|-----|-----|-------|-------|----|-----|-------|-------|-------|-----|----|
| لا        | بُّ | بَّ | بّ    | بْ    | بُ | بِي | بَا   | بُ    | بِ    | بَا |    |
| lām 'alif |     |     | šadda | sukūn |    |     | damma | kasra | fatha |     |    |
| lā        | bbu | bbi | bba   | bb    | b  | bū  | bī    | bā    | bu    | bi  | ba |

What are diacritics?

- 1 Short vowels: *fatha*, *kasra* and *damma*
- 2 Non-vowel mark: *sukun*
- 3 Double consonant mark: *shadda*

# Language pair

Arabic: diacritics

|           |     |     |       |       |    |     |       |       |       |     |    |
|-----------|-----|-----|-------|-------|----|-----|-------|-------|-------|-----|----|
| لا        | بُ  | بِّ | بَّ   | بْ    | بُ | بِي | بَا   | بُ    | بِ    | بَا |    |
| lām 'alif |     |     | šadda | sukūn |    |     | damma | kasra | fatha |     |    |
| lā        | bbu | bbi | bba   | bb    | b  | bū  | bī    | bā    | bu    | bi  | ba |

What are diacritics?

- 1 Short vowels: *fatha*, *kasra* and *damma*
- 2 Non-vowel mark: *sukun*
- 3 Double consonant mark: *shadda*

👉 **But** diacritics are not usually seen in written texts:  
MT corpora are non-vocalized and non-diacritized.

# Language pair

Arabic: diacritics and ambiguity

The absence of diacritics increases the ambiguity

ktb

katab



to write

katib



writer

kitab

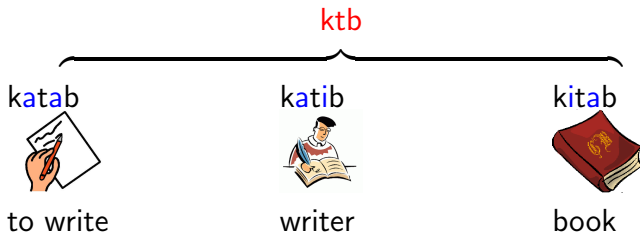


book

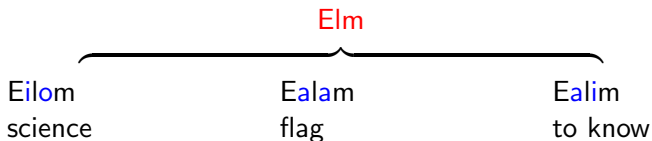
# Language pair

Arabic: diacritics and ambiguity

The absence of diacritics increases the ambiguity



Sometimes with completely different meanings:





# Language pair

Arabic: agglutination and segmentation

1 Arabic token = proclitics + affixes + root + enclitics

Example: *وبه سناتهم* (*wbHsnAthm* transliteration)  
and by their virtues

| enclitic | affix | stem       | proclitics |       |
|----------|-------|------------|------------|-------|
| hm       | At    | <b>Hsn</b> | b          | w     |
| (their)  | (s)   | (virtue)   | (by)       | (and) |

# Language pair

Arabic: agglutination and segmentation

1 Arabic token = proclitics + affixes + root + enclitics

Example: *وبه سناتهم* (*wbHsnAthm* transliteration)  
and by their virtues

| enclitic | affix | stem       | proclitics |       |
|----------|-------|------------|------------|-------|
| hm       | At    | <b>Hsn</b> | b          | w     |
| (their)  | (s)   | (virtue)   | (by)       | (and) |

☞ Enclitics: pronouns and possessives

# Language pair

Arabic: agglutination and segmentation

1 Arabic token = proclitics + affixes + root + enclitics

Example: *وبه سناتهم* (*wbHsnAthm* transliteration)  
and by their virtues

| enclitic | affix | stem       | proclitics |       |
|----------|-------|------------|------------|-------|
| hm       | At    | <b>Hsn</b> | b          | w     |
| (their)  | (s)   | (virtue)   | (by)       | (and) |

☞ Affixes: tense, genus and number marks

# Language pair

Arabic: agglutination and segmentation

1 Arabic token = proclitics + affixes + root + enclitics

Example: *وبه سناتهم* (*wbHsnAthm* transliteration)  
and by their virtues

| enclitic | affix | stem       | proclitics |          |
|----------|-------|------------|------------|----------|
| hm       | At    | <b>Hsn</b> | <b>b</b>   | <b>w</b> |
| (their)  | (s)   | (virtue)   | (by)       | (and)    |

☞ Proclitics: prepositions, conjunctions and determiners

# *Arabic-to-English translation system*

- 1 Introduction
- 2 System Design
  - Corpora
  - Pre-process
  - SMT system
- 3 Experiments and evaluation
- 4 Conclusions

- 1 News domain
  - Compilation of corpora supplied by LDC for the 2008 NIST Machine Translation Open Evaluation

- 1 News domain
  - Compilation of corpora supplied by LDC for the 2008 NIST Machine Translation Open Evaluation

| Corpus                              | Lines          | Arabic tokens    | English tokens   |
|-------------------------------------|----------------|------------------|------------------|
| Arabic English Parallel News Part 1 | 61,000         | 2,179,289        | 2,273,021        |
| Arabic News Translation Text Part 1 | 18,000         | 532,771          | 602,262          |
| Arabic Treebank English Translation | 23,800         | 660,821          | 739,695          |
| eTIRR Arabic English News Text      | 4,000          | 97,882           | 98,655           |
| Multiple-Translation Arabic         | 15,533         | 434,465          | 507,617          |
| TIDES MT2004 Arabic evaluation data | 1,329          | 40,667           | 47,324           |
| <b>Total:</b>                       | <b>123,662</b> | <b>3,945,895</b> | <b>4,262,740</b> |

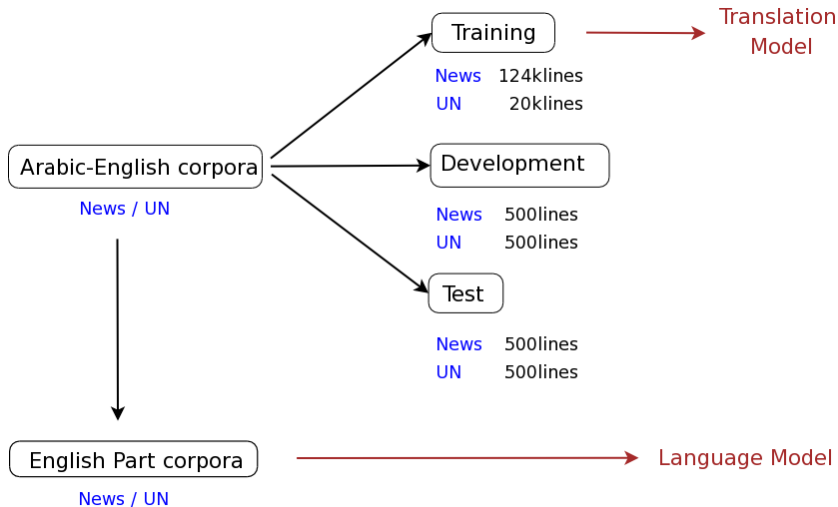
- 1 **News** domain
  - Compilation of corpora supplied by LDC for the 2008 NIST Machine Translation Open Evaluation
  - 123,662 lines, 3.9M Arabic tokens, 4.2M English tokens
- 2 **United Nations** transcriptions



- 1 **News** domain
  - Compilation of corpora supplied by LDC for the 2008 NIST Machine Translation Open Evaluation
  - 123,662 lines, 3.9M Arabic tokens, 4.2M English tokens
- 2 **United Nations** transcriptions
  - Transcriptions from 1993 to 2002
  - Whole corpus: 3,686,372 lines
  - Training set: 20,000 lines, 805K Arabic tokens, 642K English tokens

# Corpora

## Parallel corpora



# Linguistic processing

## Arabic

*First step:*

### Buckwalter transliteration

- One to one correspondence between the Arabic glyphs and UTF-8
- Replace XML characters

|   |   |   |
|---|---|---|
| ء | ذ | ل |
| أ | ر | م |
| أ | ز | ن |
| ؤ | س | ه |
| إ | ش | و |
| ئ | ص | ي |
| ا | ض | ي |
| ب | ط | ف |
| ة | ظ | ن |
| ت | ع | ك |
| ث | غ | ا |
| ج | ـ | u |
| ح | ف | i |
| خ | ق | ~ |
| د | ك | o |

### Original text:

وتأتي دول فرنسا وبريطانيا وإيطاليا وألمانيا وأيرلندا وأسبانيا ولو كسمبورج في المقدمة وبخلاف الشركات الأوروبية فقد وصل حجم رءوس الأموال المصدرة لل شركات العاملة في مصر حتي ديسمبر ٢٠٠٠ الي ١٢٦ مليار ج نيه لعدد ١٠ آلاف شركة استثمارية مؤسسه وفقا لقانون الاستثمار ..

### Transliterated text:

wt>ty dwl frnsA wbryTAnyA wAyTAlYA w>lmAnyA w>yrlndA  
w>sbAnyA wlwksmbwrj fy Almqdmw bxlAf Al\$rkAt Al>wrwbyp  
fqd wSl HjM r'ws Al>mwAl AlmSdrp ll\$rkAt AlEAmly fy mSr  
Hty dysbmbR 2000 Aly 126 mlyAr jnyh lEdd 10 |lAf \$rkp  
AstvmAryp m&ssp wfqA lqAnwn AlAstvmAr ..

*Second step:*

### Tokenization

- Word segmentation in proclitics, stems+affixes, and enclitics
- Separates punctuation
- ASVMTools (Diab 2004) –except *Al-* determiner–

Example:

```
w t>ty dwl frnsA w bryTAnyA wAyTAlYA w >lmAnyA w >yrlndA w  
>sbAnyA w lwksmbwrj fy AlmQdmp w b xLAf Al$rkAt Al>wrwbyp  
f qd wSl Hjm r'ws Al>mwAl AlmSdrp l Al$rkAt AlEAmIp fy  
mSr Hty dysbmr 2000 Aly 126 mlyAr jnyh l Edd 10 |lAf $rkp  
AstvmAryp m&ssp wfqA l qAnwn AlAstvmAr
```

First + Second step  $\implies$  Linguistic tokens

*Third step:*

Annotate with Part-of-Speech and Chunk

- ASVMTools (Diab 2004)
- PoS: 24 tags (noun, adjective, verb...)
- Chunk: IOB tagging scheme (Inside-Outside-Beginning)
- Final text: *word|lemma|PoS|chunk*

### Example:

w|w|CC|O tOty|tOty|VBP|B-VP dwl|dwl|NN|B-NP frnsA|frnsA|NNP|B-NP  
w|w|CC|O bryTAnyA|bryTAnyA|NNP|B-NP wAyTAlYA|wAyTAlYA|JJ|I-NP w|w|CC|O  
OlmAnyA|OlmAnyA|NNP|B-NP w|w|CC|O OyrIndA|OyrIndA|NNP|B-NP w|w|CC|O  
OsbAnyA|OsbAnyA|NNP|B-NP w|w|CC|O lwksmbwrj|lwksmbwrj|NNP|B-NP  
fy|fy|IN|B-PP Almqdmp|Almqdmp|NN|B-NP w|w|CC|B-PP b|b|IN|B-PP  
xLAf|xLAf|NN|B-NP Al\$rkAt|Al\$rkAt|NNS|B-NP AlOwrwbyp|AlOwrwbyp|JJ|I-NP  
f|f|CC|B-ADVP qd|qd|RP|B-PRT wSl|wSl|VBD|B-VP HjM|HjM|NN|B-NP  
r'ws|r'ws|NN|B-NP AlOmWAl|AlOmWAl|NN|B-NP AlmSdrp|AlmSdrp|JJ|B-ADJP  
l|l|IN|B-PP Al\$rkAt|Al\$rkAt|NNS|B-NP AlEAmlp|AlEAmlp|JJ|I-NP  
fy|fy|IN|B-PP mSr|mSr|NNP|B-NP Hty|Hty|IN|B-PP dysbmr|dysbmr|NN|B-NP  
2000|2000|CD|B-NP Aly|Aly|IN|B-PP 126|126|CD|B-NP mlyAr|mlyAr|NN|I-NP  
jnyh|jnyh|NN|I-NP l|l|IN|B-PP Edd|Edd|NN|B-NP 10|10|CD|B-NP  
LlAf|LlAf|NN|I-NP \$rkp|\$rkp|NN|I-NP AstvmAryp|AstvmAryp|JJ|I-NP  
mWssp|mWssp|NN|B-NP wfqA|wfqA|NN|B-NP l|l|IN|B-PP qAnwn|qAnwn|NN|B-NP  
AlAstvmAr|AlAstvmAr|NN|B-NP .|. |PUNC|O .|. |PUNC|O

*First step:*

Lowercase & Tokenization

*Second step:*

Part-of-Speech

- SVMTool (Giménez & Màrquez 2004)
- 36 tags (noun, adjective, verb...)

*Third step:*

Lemmatization

- Table (word,PoS) → lemma (185,201 entries)



*Fourth step:*

### Chunking

- Yamcha (Kudo 2003)
- IOB tagging scheme (Inside-Outside-Beginning)
- Final text: *word|lemma|PoS|chunk*

# Linguistic processing

## English

### Example:

france|france|NN|B-NP ,|,|,|I-NP britain|britain|NN|I-NP ,|,|,|O  
italy|italy|RB|B-ADVP ,|,|,|O germany|germany|NN|B-NP ,|,|,|O  
ireland|ireland|NN|B-NP ,|,|,|O spain|spain|NN|B-NP ,|,|,|O and|and|CC|O  
luxembourg|luxembourg|NN|B-NP came|come|VBD|B-VP first.|first.|RB|B-ADVP  
a|a|DT|B-NP part|part|NN|I-NP from|from|IN|B-PP the|the|DT|B-NP  
european|european|JJ|I-NP companies|company|NNS|I-NP ,|,|,|O  
the|the|DT|B-NP issued|issue|VBN|I-NP capital|capital|NN|I-NP  
of|of|IN|B-PP companies|company|NNS|B-NP operating|operate|VBG|B-VP  
in|in|IN|B-PP egypt|egypt|NN|B-NP reached|reach|VBN|B-VP  
1e126|1e126|NN|B-NP billion|billion|CD|I-NP up|up|RP|B-ADVP  
till|till|IN|B-PP december|december|NN|B-NP 2000.|2000.|CD|I-NP  
such|such|JJ|I-NP capital|capital|NN|I-NP is|be|VBZ|B-VP of|of|IN|B-PP  
10,000|10,000|CD|B-NP investment|investment|NN|I-NP  
companies|company|NNS|I-NP set|set|VBN|B-VP up|up|RP|B-PRT  
under|under|IN|B-PP the|the|DT|B-NP investment|investment|NN|I-NP  
law|law|NN|I-NP .|.|.|O

# Statistical Machine Translation System

## Building the system

- ☞ Language model
  - ☞ 5-gram Language Model, interpolated Kneser-Ney discounting
  - ☞ SRILM Toolkit (Stolcke 2002)
- Translation model
  - ▶ Alignments: GIZA++ Toolkit (Och & Ney 2003)
  - ▶ Translation tables: Moses package (Koehn et al. 2006)  
MLT package (Giménez & Màrquez)
- Decoder
  - ▶ Moses decoder (Koehn et al. 2006)
- Weights optimization
  - ▶ MERT (reference score BLEU)

# Statistical Machine Translation System

## Building the system

- Language model

- ▶ 5-gram Language Model, interpolated Kneser-Ney discounting
- ▶ SRILM Toolkit (Stolcke 2002)

- ☞ Translation model

- ☞ Alignments: GIZA++ Toolkit (Och & Ney 2003)
- ☞ Translation tables: Moses package (Koehn et al. 2006)  
MLT package (Giménez & Màrquez)

- Decoder

- ▶ Moses decoder (Koehn et al. 2006)

- Weights optimization

- ▶ MERT (reference score BLEU)

# Statistical Machine Translation System

## Building the system

- Language model
  - ▶ 5-gram Language Model, interpolated Kneser-Ney discounting
  - ▶ SRILM Toolkit (Stolcke 2002)
- Translation model
  - ▶ Alignments: GIZA++ Toolkit (Och & Ney 2003)
  - ▶ Translation tables: Moses package (Koehn et al. 2006)  
MLT package (Giménez & Màrquez)
- ☞ Decoder
  - ☞ Moses decoder (Koehn et al. 2006)
- Weights optimization
  - ▶ MERT (reference score BLEU)

# Statistical Machine Translation System

## Building the system

- Language model
  - ▶ 5-gram Language Model, interpolated Kneser-Ney discounting
  - ▶ SRILM Toolkit (Stolcke 2002)
- Translation model
  - ▶ Alignments: GIZA++ Toolkit (Och & Ney 2003)
  - ▶ Translation tables: Moses package (Koehn et al. 2006)  
MLT package (Giménez & Màrquez)
- Decoder
  - ▶ Moses decoder (Koehn et al. 2006)
- ☞ Weights optimization
  - ☞ MERT (reference score BLEU)

- 1 Introduction
- 2 System Design
- 3 Experiments and evaluation**
  - Word segmentation in Arabic
  - SMT System with Linguistic Information
  - SMT System with Discriminative Phrase Selection
- 4 Conclusions

# Experiments and evaluation

## Word segmentation in Arabic

Segmentation in the *News* compilation:

|Al\$rkAt vs. | Al\$rkAt vs. | Al \$rkAt

|                   | lines   | tokens    | toks/line |
|-------------------|---------|-----------|-----------|
| punct.            | 124,154 | 3,402,824 | 27.4      |
| punct.+clitics    | 123,662 | 3,939,726 | 31.8      |
| punct.+clitics+Al | 123,498 | 4,718,933 | 38.2      |
| English           | 123,662 | 4,262,740 | 34.5      |



# Experiments and evaluation

## Word segmentation in Arabic

Evaluation: BLEU ( $n$ -gram based metric)

For the three levels of segmentation:

|                   | Arabic→English |       | English→Arabic |       |
|-------------------|----------------|-------|----------------|-------|
|                   | dev            | test  | dev            | test  |
| punct.            | 25.76          | 23.46 | 23.50          | 16.17 |
| punct.+clitics    | 26.25          | 23.81 | 26.54          | 19.67 |
| punct.+clitics+AI | 25.28          | 23.21 | 32.46          | 26.68 |

# Experiments and evaluation

## Word segmentation in Arabic

Evaluation: BLEU ( $n$ -gram based metric)

For the three levels of segmentation:

|                   | Arabic→English |       | English→Arabic |       |
|-------------------|----------------|-------|----------------|-------|
|                   | dev            | test  | dev            | test  |
| punct.            | 25.76          | 23.46 | 23.50          | 16.17 |
| punct.+clitics    | 26.25          | 23.81 | 26.54          | 19.67 |
| punct.+clitics+AI | 25.28          | 23.21 | 32.46          | 26.68 |

☞ Best results when similar sentence lengths for both languages

# SMT system

Including linguistic information

Including linguistic information into a standard SMT

# SMT system

Including linguistic information

Including linguistic information into a standard SMT

Methods:

- Direct concatenation of the information
- Combination of two translation models
- One translation model with two factors

# SMT system

Including linguistic information

Including linguistic information into a standard SMT

Methods:

- Direct concatenation of the information
- Combination of two translation models
- One translation model with two factors

Best BLEU results: **word&lemma**

|              | Arabic→English |              | English→Arabic |              |
|--------------|----------------|--------------|----------------|--------------|
|              | dev            | test         | dev            | test         |
| w (baseline) | 24.70          | 23.82        | 26.83          | 22.85        |
| wl           | 24.74          | <b>24.28</b> | 26.95          | <b>23.34</b> |

# SMT system with discriminative phrase selection

General idea

## Word Sense Disambiguation (WSD)

Identify the correct sense of a word given a sentence

# SMT system with discriminative phrase selection

General idea

Word Sense Disambiguation (WSD)

Identify the correct sense of a word given a sentence



Different phrase **translations**  $\equiv$  Different phrase **senses**

# SMT system with discriminative phrase selection

## General idea

Word Sense Disambiguation (WSD)

Identify the correct sense of a word given a sentence



Different phrase **translations**  $\equiv$  Different phrase **senses**



Discriminative Phrase Translation (DPT)



# SMT system with discriminative phrase selection

## The method

Discriminative phrase selection:

- ☞ Phrase selection is treated as a classification problem
  - We use SVMs to solve the multiclass classification problem

# SMT system with discriminative phrase selection

## The method

Discriminative phrase selection:

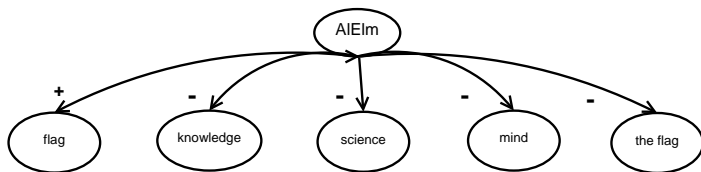
- Phrase selection is treated as a classification problem
- 👉 We use SVMs to solve the multiclass classification problem

# SMT system with discriminative phrase selection

## The method

Discriminative phrase selection:

- Phrase selection is treated as a classification problem
- We use SVMs to solve the multiclass classification problem
- ☞ Every possible translation is a class → one-vs-all classification:



# SMT system with discriminative phrase selection

## The method

SVMs allow to use **context** and **linguistic information**

Features set for the SVMs include:

- Source phrase features
  - ▶ Word, PoS, coarse PoS and chunk  $n$ -grams
- Source sentence features
  - ▶ Word, PoS, coarse PoS, chunk  $n$ -grams and bag-of-words

For the previous example, AIEIm:

# SMT system with discriminative phrase selection

## The method

Since the phrase is a word, phrase features are 1-grams:

### Sentence:

---

w tAbE mr\$d AllxwAn " In AIEIm AlmTlwb fy dyn nA hw kl Elm  
nAfE tbqY l AlnAs vmrt h , swA' kAn ElmAF \$rEyAF  
Ow ElmAF tjrybyAF .

### Phrase features:

---

|                       |       |
|-----------------------|-------|
| word $n$ -grams       | AIEIm |
| PoS $n$ -grams        | NN    |
| coarse PoS $n$ -grams | N     |
| chunk $n$ -grams      | B-NP  |

# SMT system with discriminative phrase selection

## The method

And the context of the whole sentence:

### Sentence features:

---

|                               |   |
|-------------------------------|---|
| word                          | (AlmTlwb) <sub>1</sub> , (fy) <sub>2</sub> , (dyn) <sub>3</sub> , (nA) <sub>4</sub> , (hw) <sub>5</sub> , <i>n</i> -grams (" In) <sub>-2</sub> , (AllxwAn) <sub>-3</sub> , (mr\$d) <sub>-4</sub> , (tAbE) <sub>-5</sub> , (AlmTlwb fy) <sub>1</sub> , (fy dyn) <sub>2</sub> , (dyn nA) <sub>3</sub> , (nA hw) <sub>4</sub> , (In AlmTlwb) <sub>-1</sub> , (AllxwAn ") <sub>-3</sub> , (mr\$d AllxwAn) <sub>-4</sub> , (tAbEmr\$d) <sub>-5</sub> , (AlmTlwb fy dyn) <sub>1</sub> , (fy dyn nA) <sub>2</sub> , (dyn nA hw) <sub>3</sub> , (In AlmTlwb fy) <sub>-1</sub> , (" In AlmTlwb) <sub>-2</sub> , (AllxwAn " In) <sub>-3</sub> , (mr\$d AllxwAn ") <sub>-4</sub> , (tAbE mr\$d AllxwAn) <sub>-5</sub>  |
| PoS                           | (JJ) <sub>1</sub> , (IN) <sub>2</sub> , (NN) <sub>3</sub> , (PRP\$) <sub>4</sub> , (PRP) <sub>5</sub> , <i>n</i> -grams (PUNC IN) <sub>-2</sub> , (NN) <sub>-3</sub> , (NN) <sub>-4</sub> , (VBD) <sub>-5</sub> , (JJ IN) <sub>1</sub> , (IN NN) <sub>2</sub> , (NN PRP\$) <sub>3</sub> , (PRP\$ PRP) <sub>4</sub> , (IN JJ) <sub>-1</sub> , (NN PUNC) <sub>-3</sub> , (NN NN) <sub>-4</sub> , (VBD NN) <sub>-5</sub> , (JJ IN NN) <sub>1</sub> , (IN NN PRP\$) <sub>2</sub> , (NN PRP\$ PRP) <sub>3</sub> , (IN JJ IN) <sub>-1</sub> , (PUNC IN JJ) <sub>-2</sub> , (NN PUNC IN) <sub>-3</sub> , (NN NN PUNC) <sub>-4</sub> , (VBD NN NN) <sub>-5</sub> , (J) <sub>1</sub> , (I) <sub>2</sub> , (N) <sub>3</sub> , (P) <sub>4</sub> , (P) <sub>5</sub> , (P I) <sub>-2</sub> , (N) <sub>-3</sub> , (N) <sub>-4</sub> , (V) <sub>-5</sub> |
| coarse PoS<br><i>n</i> -grams | (J I) <sub>1</sub> , (I N) <sub>2</sub> , (N P) <sub>3</sub> , (P P) <sub>4</sub> , (I J) <sub>-1</sub> , (N P) <sub>-3</sub> , (N N) <sub>-4</sub> , (V N) <sub>-5</sub> , (J I N) <sub>1</sub> , (I N P) <sub>2</sub> , (N P P) <sub>3</sub> , (I J I) <sub>-1</sub> , (P I J) <sub>-2</sub> , (N P I) <sub>-3</sub> , (N N P) <sub>-4</sub> , (V N N) <sub>-5</sub>  |
| chunk<br><i>n</i> -grams      | (I-NP) <sub>1</sub> , (B-PP) <sub>2</sub> , (B-NP) <sub>3</sub> , (I-NP) <sub>4</sub> , (B-NP) <sub>5</sub> , (O B-SBAR) <sub>-2</sub> , (B-NP) <sub>-3</sub> , (B-NP) <sub>-4</sub> , (B-VP) <sub>-5</sub> , (I-NP B-PP) <sub>1</sub> , (B-PP B-NP) <sub>2</sub> , (B-NP I-NP) <sub>3</sub> , (I-NP B-NP) <sub>4</sub> , (B-SBAR I-NP) <sub>-1</sub> , (B-NP O) <sub>-3</sub> , (B-NP B-NP) <sub>-4</sub> , (B-VP B-NP) <sub>-5</sub> , (I-NP B-PP B-NP) <sub>1</sub> , (B-PP B-NP I-NP) <sub>2</sub> , (B-NP I-NP B-NP) <sub>3</sub> , (B-SBAR I-NP B-PP) <sub>-1</sub> , (O B-SBAR I-NP) <sub>-2</sub> , (B-NP O B-SBAR) <sub>-3</sub> , (B-NP B-NP O) <sub>-4</sub> , (B-VP B-NP B-NP) <sub>-5</sub>  |
| bag-of-words                  | left: AllxwAn, mr\$d, tAbE<br>right: \$rEyAF, AlmTlwb, AlnAs, Elm, ElmAF, dyn, kAn, kl, nAfE, swA', tbqY, tjrybyAF, vmrt  |

# SMT system with discriminative phrase selection

## The system

Estimation of the discriminative phrase translation model and integration into the SMT system:

- ☞ Training linear SVMs (SVM<sup>light</sup>, Joachims 1999) for every translation of every phrase
- Convert SVM score into probability via a *softmax* function
- Include this probability in the translation model within a Log-linear model

# SMT system with discriminative phrase selection

## The system

Estimation of the discriminative phrase translation model and integration into the SMT system:

- Training linear SVMs ( $SVM^{light}$ , Joachims 1999) for every translation of every phrase
- ☞ Convert SVM score into probability via a *softmax* function
- Include this probability in the translation model within a Log-linear model



# SMT system with discriminative phrase selection

## The system

Estimation of the discriminative phrase translation model and integration into the SMT system:

- Training linear SVMs ( $SVM^{light}$ , Joachims 1999) for every translation of every phrase
- Convert SVM score into probability via a *softmax* function
- 👉 Include this probability in the translation model within a Log-linear model

# SMT system with discriminative phrase selection

## Discriminative phrase translation

### Phrase translation task

Improvement in accuracy wrt. the most frequent translation, MFT

| Occurrences   | #     | Acc.DPT (%) | Acc.MFT (%) |
|---------------|-------|-------------|-------------|
| 100-500       | 4,310 | 66.5        | 58.7        |
| 501-1,000     | 565   | 68.8        | 62.3        |
| 1,001-5,000   | 393   | 73.0        | 66.7        |
| 5,001-10,000  | 27    | 79.5        | 72.2        |
| 10,001-50,000 | 19    | 74.8        | 66.6        |
| > 50,000      | 7     | 80.7        | 76.2        |
| Total:        | 5,321 | <b>67.3</b> | <b>59.8</b> |

# SMT system with discriminative phrase selection

Integration into the SMT system

## Full translation task

Hyv<sub>28</sub> tm<sub>22</sub> AHrAq AIEIm<sub>1</sub> AldnmArky .1128

Translation table example:

| $f_i$              | $e_j$       | $P_{DPT}(e f)$ | $P_{MLE}(e f)$ |
|--------------------|-------------|----------------|----------------|
| AIEIm <sub>1</sub> | flag        | 0.1986         | 0.3241         |
| AIEIm <sub>1</sub> | the         | 0.0419         | 0.0207         |
| AIEIm <sub>1</sub> | mind        | 0.0401         | 0.0620         |
| AIEIm <sub>1</sub> | the flag    | 0.0397         | 0.0414         |
| AIEIm <sub>1</sub> | flag during | 0.0394         | 0.0138         |
| AIEIm <sub>1</sub> | knowledge   | 0.0392         | 0.1103         |
| AIEIm <sub>1</sub> | flag caused | 0.0387         | 0.0138         |
| AIEIm <sub>1</sub> | science     | 0.0377         | 0.1793         |
| AIEIm <sub>1</sub> | education   | 0.0377         | 0.0138         |
| AIEIm <sub>1</sub> | in mind     | 0.0371         | 0.0138         |

# SMT system with discriminative phrase selection

Integration into the SMT system

## Full translation task

Hyv<sub>28</sub> tm<sub>22</sub> AHrAq **AlElm<sub>1</sub>** AldnmArky .1128

Translation table example:

| $f_i$              | $e_j$       | $P_{DPT}(e f)$ | $P_{MLE}(e f)$ |
|--------------------|-------------|----------------|----------------|
| AlElm <sub>1</sub> | flag        | <b>0.1986</b>  | <b>0.3241</b>  |
| AlElm <sub>1</sub> | the         | 0.0419         | 0.0207         |
| AlElm <sub>1</sub> | mind        | 0.0401         | 0.0620         |
| AlElm <sub>1</sub> | the flag    | 0.0397         | 0.0414         |
| AlElm <sub>1</sub> | flag during | 0.0394         | 0.0138         |
| AlElm <sub>1</sub> | knowledge   | 0.0392         | <b>0.1103</b>  |
| AlElm <sub>1</sub> | flag caused | 0.0387         | 0.0138         |
| AlElm <sub>1</sub> | science     | 0.0377         | <b>0.1793</b>  |
| AlElm <sub>1</sub> | education   | 0.0377         | 0.0138         |
| AlElm <sub>1</sub> | in mind     | 0.0371         | 0.0138         |

# SMT system with discriminative phrase selection

## Evaluation

Three systems:

SMT  
standard

DPT  
replace MLE

DPT<sup>+</sup>  
add

# SMT system with discriminative phrase selection

## Evaluation

Three systems:

SMT  
standard

DPT  
replace MLE

DPT<sup>+</sup>  
add

Results:

|      | SMT   | DPT   | DPT <sup>+</sup> |
|------|-------|-------|------------------|
| BLEU | 23.88 | 23.87 | 23.96            |

👉 Non-significative improvements obtained

# SMT system with discriminative phrase selection

## Evaluation

Three systems:

SMT  
standard

DPT  
replace MLE

DPT<sup>+</sup>  
add

Results:

|      | SMT   | DPT   | DPT <sup>+</sup> |
|------|-------|-------|------------------|
| BLEU | 23.88 | 23.87 | 23.96            |

- Non-significative improvements obtained
- ☞ Coherent with other metrics, but in general, DPT better than DPT<sup>+</sup>

# *Summary and conclusions*

- 1 Introduction
- 2 System Design
- 3 Experiments and evaluation
- 4 Conclusions**
  - Summary
  - Future work



## Summary

- ☞ First approach to the Arabic-to-English translation task
  - Parallel corpora annotated with lemma, PoS and chunk
  - Clitic segmentation improves the translation performance
  - A direct inclusion of linguistic information (lemmas) slightly improves the results

## Summary

- First approach to the Arabic-to-English translation task
- 👉 Parallel corpora annotated with lemma, PoS and chunk
- Clitic segmentation improves the translation performance
- A direct inclusion of linguistic information (lemmas) slightly improves the results

## Summary

- First approach to the Arabic-to-English translation task
- Parallel corpora annotated with lemma, PoS and chunk
- 👉 Clitic segmentation improves the translation performance
- A direct inclusion of linguistic information (lemmas) slightly improves the results

## Summary

- First approach to the Arabic-to-English translation task
- Parallel corpora annotated with lemma, PoS and chunk
- Clitic segmentation improves the translation performance
- 👉 A direct inclusion of linguistic information (lemmas) slightly improves the results

## Summary

- 👉 We try to apply WSD techniques to the translation task. That allows to consider linguistic information and the context of the phrase
- Although we get an increment of 7.5% in accuracy in the phrase translation task, this is not reflected in full translation
- Necessity to improve the integration of DPT predictions into SMT system

## Summary

- We try to apply WSD techniques to the translation task. That allows to consider linguistic information and the context of the phrase
- ☞ Although we get an increment of 7.5% in accuracy in the phrase translation task, this is not reflected in full translation
- Necessity to improve the integration of DPT predictions into SMT system

## Summary

- We try to apply WSD techniques to the translation task. That allows to consider linguistic information and the context of the phrase
- Although we get an increment of 7.5% in accuracy in the phrase translation task, this is not reflected in full translation
- 👉 Necessity to improve the integration of DPT predictions into SMT system

## Future work

- ☞ Complete DPT scores in the translation table when DPT prediction is not available
- Complement DPT predictions for both directions,  $P_{\text{DPT}}(f|e)$  and  $P_{\text{DPT}}(e|f)$ , as done with MLE probabilities
- Deep into Arabic processing, especially in tokenization and lemmatization



## Future work

- Complete DPT scores in the translation table when DPT prediction is not available
- 👉 Complement DPT predictions for both directions,  $P_{\text{DPT}}(f|e)$  and  $P_{\text{DPT}}(e|f)$ , as done with MLE probabilities
- Deep into Arabic processing, especially in tokenization and lemmatization

## Future work

- Complete DPT scores in the translation table when DPT prediction is not available
- Complement DPT predictions for both directions,  $P_{\text{DPT}}(f|e)$  and  $P_{\text{DPT}}(e|f)$ , as done with MLE probabilities
- 👉 Deep into Arabic processing, especially in tokenization and lemmatization

OpenMT project of the *Ministerio de Educación y Ciencia*

Grup de Processament del Llenguatge Natural  
Universitat Politècnica de Catalunya

...i GRÀCIES PEL SUPORT!

Jesús Giménez i Lluís Màrquez