

# A multi-genre SMT system for Arabic to French

**Saša Hasan and Hermann Ney**

**LREC 2008**

**Marrakech, Morocco – May 29, 2008**

**Human Language Technology and Pattern Recognition**

**Lehrstuhl für Informatik 6**

**Computer Science Department**

**RWTH Aachen University, Germany**

# Overview

- ▶ **Project TRAMES:**  
**Traduction Automatique par des Méthodes Statistiques**
- ▶ **Goal:**  
**online system for translation of Arabic to French**
- ▶ **Development over 3-year period (2005–2007):**
  - ▷ **corpus gathering**
  - ▷ **preprocessing pipeline**
  - ▷ **phrase-based SMT module (decoder)**
  - ▷ **fine-tuning for different genres**
  - ▷ **software engineering for “real-time” performance**

# Rough processing pipeline (1)

## ▶ Data acquisition

- ▷ no parallel corpora initially available for Arabic-French
- ▷ gather data from the web (intl. organizations, news agencies, journals)
- ▷ main data resource: Official Document System of the United Nations (ODS)

## ▶ Corpus creation

- ▷ document and sentence alignment
- ▷ preprocessing: tokenization, Arabic word segmentation

## ▶ Training the models

- ▷ word alignments
- ▷ phrase extraction
- ▷ language modeling

## Rough processing pipeline (2)

- ▶ **Generation of translations (search/decoding)**
  - ▷ phrase-based decoder using log-linear combination of models
  - ▷ dynamic programming beam search
  - ▷ tune parameters on development set using MERT
- ▶ **Experiments**
  - ▷ evaluation of the system using automatic evaluation measures
  - ▷ compare translation output to a set of reference translations
    - BLEU:  $n$ -gram precision w/ brevity penalty
    - TER: string edit distance allowing for block movements

# Corpus creation

- ▶ Document alignment as is (from web structure)
- ▶ Sentence alignment using sentence-length model and refinements from IBM model 1 probabilities
- ▶ Preprocessing:
  - ▷ tokenization and categorization for numbers, months and URLs
  - ▷ text normalization: remove diacritics
  - ▷ word segmentation:
    - prefix and suffix splitting based on finite-state automaton
  - ▷ example:

المدرسةُ	⇒	المدرسة		
<b>the school</b>				
والمدرسة	⇒	مدرسة	ال	و
		<b>school</b>	<b>the</b>	<b>and</b>
مدرستهم	⇒	هم	مدرسة	
		<b>their</b>	<b>school</b>	

# Corpus statistics

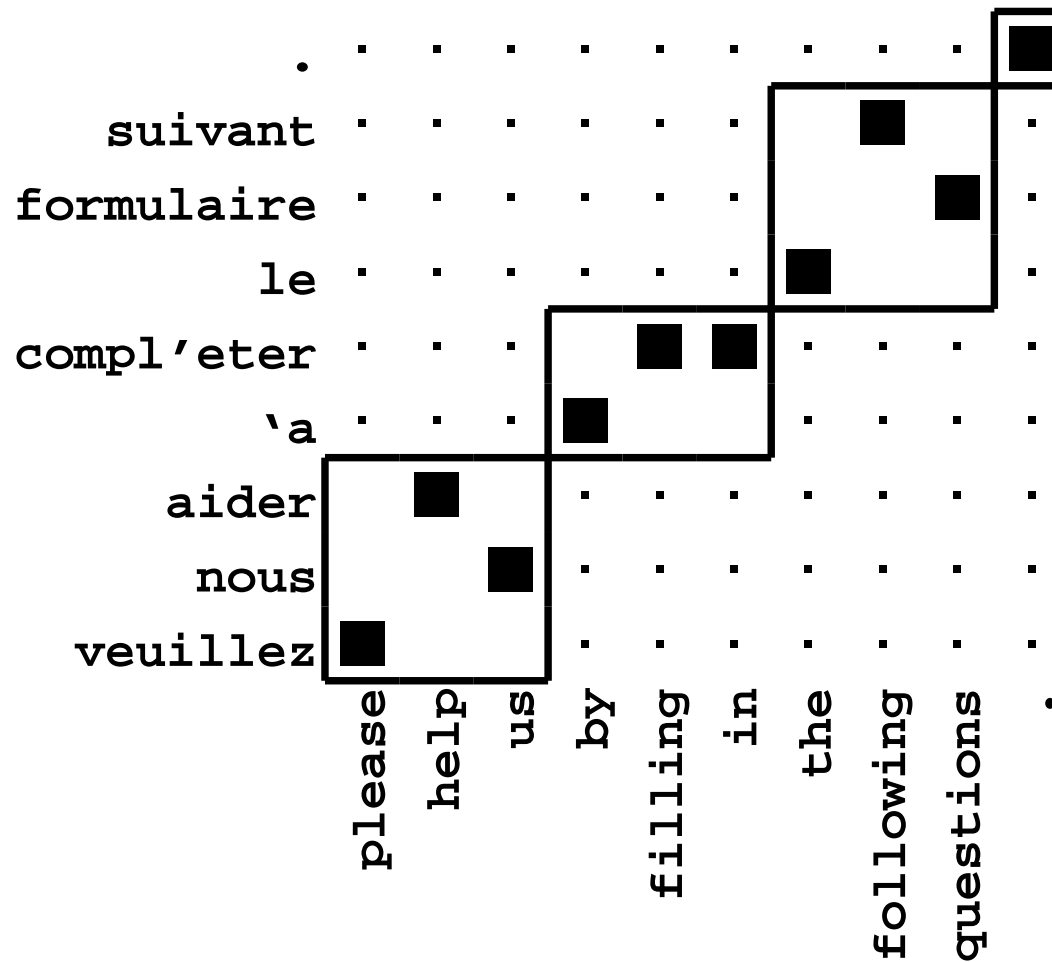
Corpus extracted from UN documents / Amnesty Int. / Le Monde Diplomatique:

	2005 system		2007 system	
	Arabic	French	Arabic	French
<b>Doc. pairs</b>	<b>62K</b>		<b>74K</b>	
<b>Sent. pairs</b>	<b>4.7M</b>		<b>6.6M</b>	
<b>Run. words</b>	<b>108.1M</b>	<b>104.8M</b>	<b>151.3M</b>	<b>180.2M</b>
<b>Vocabulary</b>	<b>245K</b>	<b>288K</b>	<b>427K</b>	<b>301K</b>

► **Important data update from BN radio and TV transcripts:**

- ▷ **Orient, Qatar, BBC, Alarabiya, Aljazeera, Alalam**
- ▷ **250 audio documents consisting of 90 hours radio and TV broadcasts**
- ▷ **21K sentences with 585K running words of domain-specific material for the audio domain**

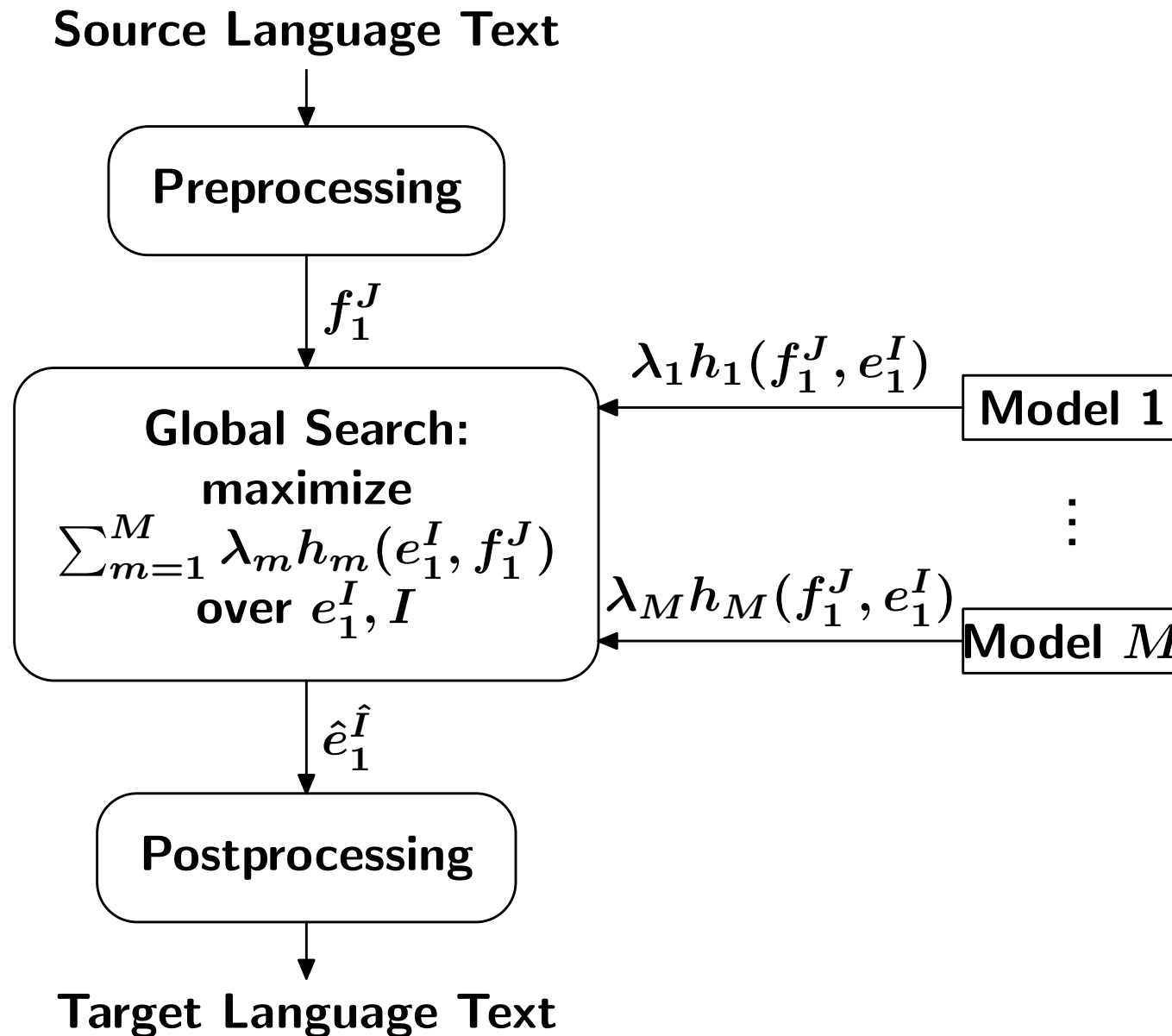
# Training and Generation (1)



Idea:

1. Segment source sentence into phrases
2. Translate each phrase
3. Concatenate these phrase translations

# Training and Generation (2)





## Evaluation: progress over time

	1st sys 2005	2nd sys 2006	+BN-LM	3rd sys 2007
<b>CESTA run2</b>	<b>40.8</b>	<b>42.9</b>	<b>43.8</b>	<b>44.8</b>
<b>Arabic BN</b>				
text setting	20.9	29.7	-	34.4
audio setting	-	34.4	37.6	41.1

- ▶ **System was tuned on held-out development sets**
- ▶ **Results shown are all on blind test sets:**
  - ▷ text domain: CESTA run2 evaluation data
  - ▷ audio domain: Arabic BN transcripts from TV/radio
- ▶ **Observations:**
  - ▷ adding BN transcripts to the system significantly boosts performance on audio
  - ▷ genre-specific tuning makes a difference

## Evaluation: comparison to Moses

	<b>BLEU [%]</b>	<b>TER [%]</b>	<b>Translation speed [words/sec]</b>
<b>CESTA run2</b>			
<b>Moses</b>	<b>42.2</b>	<b>52.25</b>	<b>14.2</b>
<b>TRAMES</b>	<b>43.4</b>	<b>51.30</b>	<b>222.0</b>
<b>Arabic BN</b>			
<b>Moses</b>	<b>39.5</b>	<b>53.37</b>	<b>18.6</b>
<b>TRAMES</b>	<b>40.0</b>	<b>52.93</b>	<b>249.3</b>

- ▶ **Freely available: open-source phrase-based decoder Moses**
- ▶ **Models / search concept similar to RWTH's decoder**
- ▶ **Fair comparison: table shows experiments for the same training data and similar pruning parameters (histogram size 200)**
- ▶ **Result: TRAMES system is up to 16 times faster with up to 250 words/sec**

## Examples: text setting

► **Arabic source sentence:**

ويتم التركيز على الوقاية من انتقال هذا المرض من الأم إلى الطفل واتخاذ نهج للنهوض بالوعي العام بين الشباب.

► **French translation, system update in 2005:**

et met l'accent sur la prévention \_\_\_ de cette maladie de la mère à l'enfant et \_\_\_  
une démarche pour la promotion de la sensibilisation du public chez les jeunes.

► **French translation, system update in 2006:**

L'accent est mis sur la prévention de la transmission \_\_\_ de la mère à l'enfant et  
une approche pour la promotion de la sensibilisation du public chez les jeunes.

► **French translation, system update in 2007:**

L'accent est mis sur la prévention de la transmission de la maladie de la mère  
à l'enfant et une approche pour promouvoir une prise de conscience parmi les  
jeunes.

► **French reference translation (1/4):**

L'accent est mis sur la prévention de la transmission de cette maladie de la mère à  
l'enfant et l'adoption de la démarche de la généralisation de la prise de conscience  
parmi les jeunes.

## Examples: audio setting

### ► Arabic source

رياض محمد رصد ردود الشارع الإيراني حيال محاكمة صدام ووافانا بالتقرير التالي.

### ► French sys1 2005

Riyad Mohammed suivi réponses la rue des UNK\_إيراني\_ pour juger Saddam  
et UNK\_وافانا\_ du rapport UNK\_التالي\_.

### ► French sys2 2006

Riad Mohamad de suivre les mesures prises par la rue iranienne par juger Saddam  
et nous a fait parvenir le rapport suivant.

### ► French sys3 2007

Riad Mohamad suivi de la réponse de la rue iranienne envers le procès de Saddam  
et nous a fait parvenir le rapport suivant.

### ► French reference translation

Riad Mohamed a scruté les réactions dans la rue iranienne au sujet du procès de  
Saddam et nous a préparé le reportage suivant.

# Conclusions

- ▶ Presented a state-of-the-art SMT system for Arabic-to-French
- ▶ Multi-genre capability:
  - ▷ newswire (text domain)
  - ▷ broadcast news transcripts (audio domain)
- ▶ Real-time translation speeds of up to 250 words/sec
- ▶ Favorable performance:
  - ▷ BLEU 44.8% on text input
  - ▷ BLEU 41.1% on audio transcripts

## Outlook:

- ▶ Further system updates with additional data
- ▶ Additional genres, e.g. web texts (e.g. weblogs, news groups)
- ▶ On-the-fly genre determination using text classification

**Thank you for your attention**

**Saša Hasan**

[hasan@cs.rwth-aachen.de](mailto:hasan@cs.rwth-aachen.de)

<http://www-i6.informatik.rwth-aachen.de/>

# Test sets

## Blind test sets:

- ▶ CESTA run2 for text
- ▶ Arabic BN for audio setting

	Text setting		Audio setting	
	Arabic	French	Arabic	French
<b>Doc. pairs</b>	<b>30</b>		<b>7</b>	
<b>Sentences</b>	<b>824</b>	<b>3 296 (4x)</b>	<b>466</b>	<b>1 864 (4x)</b>
<b>Run. words</b>	<b>22 045</b>	<b>102 087</b>	<b>16 847</b>	<b>91 557</b>
<b>Vocabulary</b>	<b>4 441</b>	<b>6 335</b>	<b>5 952</b>	<b>6 943</b>
<b>OOV rate</b>	<b>0.40%</b>	<b>-</b>	<b>1.1%</b>	<b>-</b>