

**5th SALTMIL Workshop on Minority Languages  
Genoa, 23/5/2006**

# **Statistical Machine Translation with a Small Amount of Bilingual Training Data**

**Maja Popović, Hermann Ney**

**Human Language Technology and Pattern Recognition  
Lehrstuhl für Informatik VI, Computer Science Department  
RWTH Aachen University  
D-52056 Aachen, Germany**

# Overview

- **Introduction to Statistical Machine Translation (SMT)**
- **Motivation**
- **SMT with Sparse Training Data**
- **Recent Results**
  - **Spanish-English**
  - **Serbian-English**
- **Conclusions**

# Statistical Machine Translation (SMT)

- finding a target language sequence  $\hat{e}_1^I$  given a source language sequence  $f_1^J$ :

$$\hat{e}_1^I = \operatorname{argmax}_{e_1^I} \{Pr(e_1^I | f_1^J)\}$$

- log-linear model:

$$\hat{e}_1^I = \operatorname{argmax}_{e_1^I} \left\{ \sum_{m=1}^M \lambda_m h_m(e_1^I, f_1^J) \right\}$$

- language model

$$h_m(e_1^I, f_1^J) = \log Pr(e_1^I)$$

- translation model

$$h_m(e_1^I, f_1^J) = \log Pr(f_1^J | e_1^I)$$

- ...

# Motivation

- **translation model probabilities are extracted from a bilingual parallel text - training corpus**
  - **the quality of a translation system usually depends on the size of this corpus**
    - **a large bilingual parallel corpus is often not available**
- ⇒ **strategies for exploiting limited amounts of bilingual data for statistical machine translation**

# SMT with Sparse Bilingual Training Data

- **Al-Onaizan & Germann<sup>+</sup>, 2000**
  - comparing different translation methods on a small bilingual Tetun-English corpus
  - found out that the human mind is very capable of deriving dependencies such as morphology, cognates, proper names, etc.
    - ⇒ a crucial reason for better performance of human translation
  
- **Callison-Burch & Osborne, 2003**
  - co-training method for extension of a training corpus
  - new sentence pairs are produced by multiple translation models trained on different language pairs
    - ⇒ the best improvements achieved after two or three re-training rounds

# SMT with Sparse Bilingual Training Data

- **Niessen & Ney, 2004**
  - investigating impact of the corpus size for translation from German into English
  - morpho-syntactic information and conventional dictionary are used for improving the performance
  - ⇒ acceptable translation quality even with a very small corpus
- **Matusov & Popović<sup>+</sup>, 2004**
  - translation of spontaneous speech
  - acquiring additional training data using an n-gram coverage measure
  - morpho-syntactic information
- **Popović & Ney, 2005**
  - translation of Spanish-English and Catalan-English pair
  - phrasal lexicon and morpho-syntactic information for Spanish and Catalan verbs
  - ⇒ acceptable translation quality with only one thousand task-specific sentence pairs for training

# SMT with Sparse Bilingual Training Data

- **Popović & Vilar<sup>+</sup>, 2005**
  - translation of Serbian-English pair with a very small corpus
  - morpho-syntactic information and phrasal book
  - ⇒ translation results comparable with results for other language pairs
  
- **Goldwater & McClosky, 2005**
  - translation of Czech-English pair
  - morphological information
  
- **Lopez & Resnik, 2005; Martin & Mihalcea<sup>+</sup>, 2005**
  - word alignments for languages with scarce resources
    - \* Romanian-English
    - \* Inuktitut-English
    - \* Hindi-English

# Recent Results

## - Experimental Settings -

### Spanish-English and Serbian-English language pair

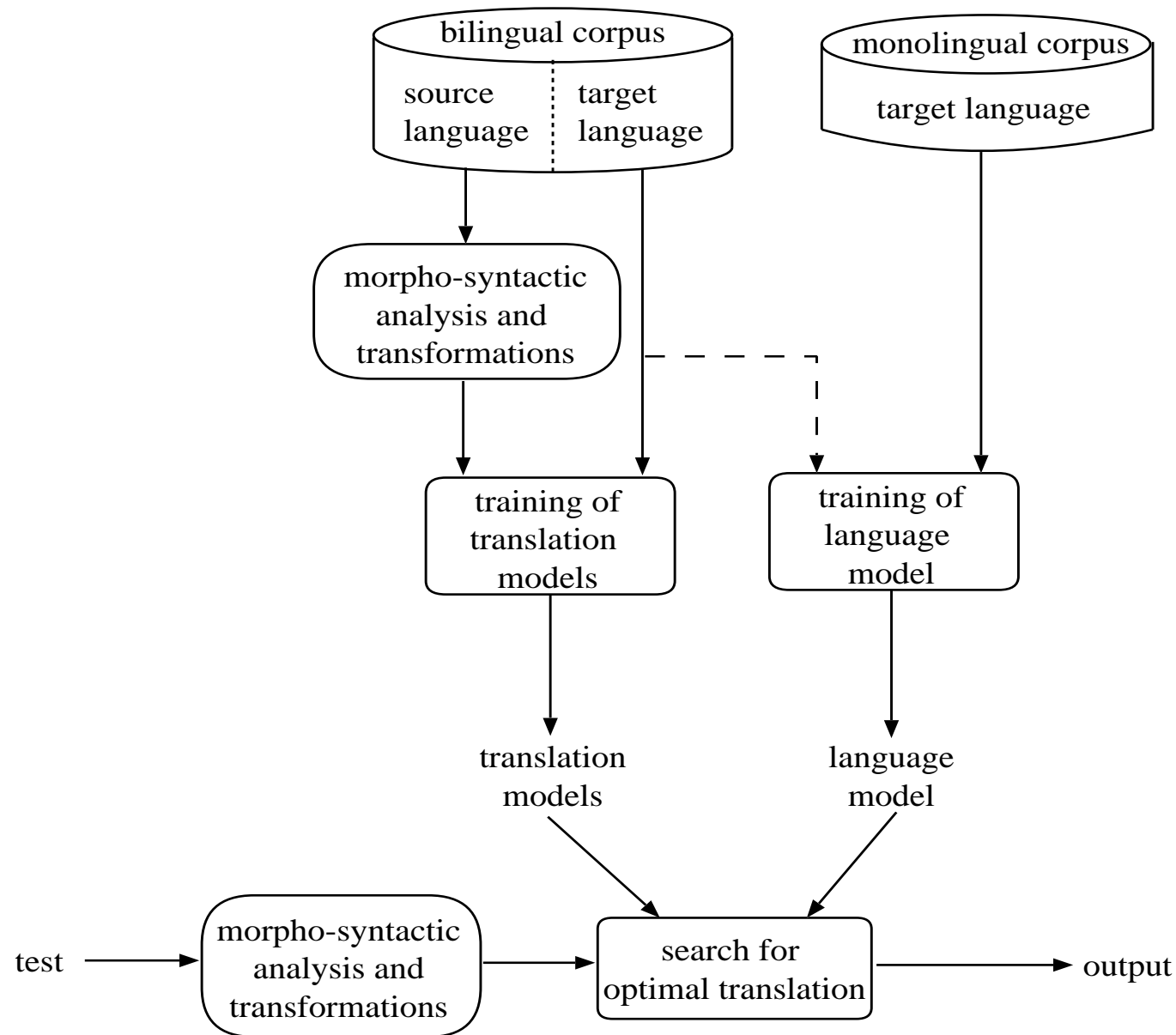
- **Spanish-English**
  - **European Parliament Plenary Sessions (EPPS) corpus**
  - **conventional dictionary**
- **Serbian-English**
  - **Assimil language course corpus**
  - **small phrasal book**
- **various sizes of bilingual training corpus**
- **appropriate morpho-syntactic transformations**
- **language model trained on the largest target language text**



# Translation System

- **state-of-the-art translation system**
- **log-linear combination of seven models:**
  - **two phrase-based models**  
**(source to target and target to source)**
  - **two single word based models at phrase level**  
**(source to target and target to source)**
  - **language model**
  - **phrase penalty and word penalty**

# Training and Test with Sparse Bilingual Resources



# Spanish↔English

## - Training Corpora -

Training		Spanish	English
<b>1.3M</b>	<b>Sentences</b>	<b>1281427</b>	
	<b>Running Words+PM</b>	<b>36578514</b>	<b>34918192</b>
	<b>Vocabulary</b>	<b>153124</b>	<b>106496</b>
	<b>Singletons [%]</b>	<b>35.2</b>	<b>36.2</b>
<b>13k</b>	<b>Sentences</b>	<b>13360</b>	
	<b>Running Words+PM</b>	<b>385198</b>	<b>366055</b>
	<b>Vocabulary</b>	<b>22425</b>	<b>16326</b>
	<b>Singletons [%]</b>	<b>47.6</b>	<b>43.7</b>
<b>1k</b>	<b>Sentences</b>	<b>1113</b>	
	<b>Running Words+PM</b>	<b>31022</b>	<b>29497</b>
	<b>Vocabulary</b>	<b>5809</b>	<b>4749</b>
	<b>Singletons [%]</b>	<b>60.8</b>	<b>55.3</b>
<b>dict.</b>	<b>Entries</b>	<b>52566</b>	
	<b>Running Words+PM</b>	<b>60964</b>	<b>62011</b>
	<b>Vocabulary</b>	<b>31126</b>	<b>30761</b>
	<b>Singletons [%]</b>	<b>67.7</b>	<b>67.4</b>

# Spanish ↔ English

## - Test Corpus -

<b>Test</b>	<b>Spanish</b>	<b>English</b>
<b>Sentences</b>	<b>840</b>	<b>1094</b>
<b>Running Words+PM</b>	<b>22774</b>	<b>26917</b>
<b>Distinct Words</b>	<b>4081</b>	<b>3958</b>
<b>OOVs (1.3M) [%]</b>	<b>0.14</b>	<b>0.25</b>
<b>OOVs (13k) [%]</b>	<b>2.8</b>	<b>2.6</b>
<b>OOVs (1k) [%]</b>	<b>10.6</b>	<b>9.4</b>
<b>OOVs (dict.) [%]</b>	<b>19.1</b>	<b>16.2</b>

# Spanish ↔ English

## - Morpho-syntactic transformations -

- local reorderings of nouns and adjectives
- replacing Spanish adjectives with their base forms

Spanish	original:	motivos <b>económicos y políticos</b>
	reordered:	<b>económicos y políticos</b> motivos
	+adjective base	<b>económico y político</b> motivos
English	original:	<b>economic and political</b> reasons
	reordered:	reasons <b>economic and political</b>

# Spanish→English

## - Translation Results -

Spanish→English		WER	PER	BLEU
<b>dict</b>	<b>baseline</b>	<b>60.4</b>	<b>49.3</b>	<b>19.4</b>
	<b>+adjective treatment</b>	<b>56.4</b>	<b>46.8</b>	<b>23.8</b>
<b>1k</b>	<b>baseline</b>	<b>52.4</b>	<b>40.7</b>	<b>30.0</b>
	<b>+dictionary</b>	<b>48.0</b>	<b>36.5</b>	<b>36.0</b>
	<b>+adjective treatment</b>	<b>44.5</b>	<b>34.8</b>	<b>40.9</b>
<b>13k</b>	<b>baseline</b>	<b>41.8</b>	<b>30.7</b>	<b>43.2</b>
	<b>+dictionary</b>	<b>40.6</b>	<b>29.6</b>	<b>46.3</b>
	<b>+adjective treatment</b>	<b>38.3</b>	<b>29.0</b>	<b>49.6</b>
<b>1.3M</b>	<b>baseline</b>	<b>34.5</b>	<b>25.5</b>	<b>54.7</b>
	<b>+reorder adjective</b>	<b>33.5</b>	<b>25.2</b>	<b>56.4</b>

- **dictionary alone might be used for multilingual information retrieval**
- **reasonable translation quality with small corpora**
  - **dictionary and morpho-syntactic information are very important**
- **1000 times larger corpus  $\Leftrightarrow$  12-25% relative decrease of error rates**

## English→Spanish - Translation Results -

English→Spanish		WER	PER	BLEU
dict	baseline	67.6	55.9	14.1
	adjective treatment	65.7	54.5	16.5
1k	baseline	60.1	47.4	23.9
	+dictionary	56.0	43.2	28.3
	adjective treatment	53.9	42.0	30.6
13k	baseline	49.6	37.4	36.2
	+dictionary	48.6	36.3	37.2
	adjective treatment	47.3	35.7	39.1
1.3M	baseline	39.7	30.6	47.8
	+reorder adjective	39.6	30.5	48.3

- similar effect as for the other translation direction
- improvements through dictionary and morpho-syntactic information are slightly smaller
  - translation into a more inflected language is more difficult
  - Spanish has a rather free word order

# Serbian↔English

## - Training Corpora -

Training		Serbian	English
<b>2.6k</b>	<b>Sentences</b>	<b>2632</b>	
	<b>Running Words+PM</b>	<b>22227</b>	<b>24808</b>
	<b>Vocabulary</b>	<b>4546</b>	<b>2645</b>
	<b>Singletons [%]</b>	<b>60.0</b>	<b>45.8</b>
<b>0.2k</b>	<b>Sentences</b>	<b>200</b>	
	<b>Running Words+PM</b>	<b>1666</b>	<b>1878</b>
	<b>Vocabulary</b>	<b>778</b>	<b>603</b>
	<b>Singletons [%]</b>	<b>79.4</b>	<b>65.5</b>
<b>phrases</b>	<b>Entries</b>	<b>351</b>	
	<b>Running Words+PM</b>	<b>617</b>	<b>730</b>
	<b>Vocabulary</b>	<b>335</b>	<b>315</b>
	<b>Singletons [%]</b>	<b>71.3</b>	<b>66.3</b>



# Serbian ↔ English

## - Test Corpus -

Test	Serbian	English
<b>Sentences</b>	<b>260</b>	
<b>Running Words+PM</b>	<b>2100</b>	<b>2336</b>
<b>Distinct Words</b>	<b>891</b>	<b>674</b>
<b>OOVs (2.6k) [%]</b>	<b>11.7</b>	<b>4.9</b>
<b>OOVs (0.2k) [%]</b>	<b>35.2</b>	<b>21.8</b>

# Serbian ↔ English

## - Morpho-syntactic transformations -

- converting Serbian words into base forms

mali		small ⇒	mali		small
mala					
malog					
malu					
malom					
...					

- additional treatment of Serbian verbs

idemo	⇒	PL1 ići		we go
idem	⇒	SG1 ići		I go

- removing English articles

when I have **the** flu, I keep **a** supply of paper handkerchiefs.



when I have flu, I keep supply of paper handkerchiefs.

# Serbian↔English

## - Translation Results -

Serbian→English		WER	PER	BLEU
0.2k	baseline	65.5	60.8	8.3
	+phrases	65.0	59.8	10.3
	+base forms	59.2	54.8	13.9
	+verb POS+neg	60.0	52.6	14.8
2.6k	baseline	44.5	37.9	32.1
	+base forms	42.9	37.4	35.4
	+verb POS+neg	41.9	34.7	34.6

- results for extremely small corpus comparable with results for a dictionary
- phrases are helpful to some extent
- morphological information is very important
- acceptable performance with less than three thousand sentence pairs

# English ↔ Serbian

## - Translation Results -

English → Serbian		WER	PER	BLEU
0.2k	baseline	73.4	68.4	6.8
	+phrases	71.9	67.5	9.3
	+remove article	66.7	62.2	9.4
2.6k	baseline	51.8	45.8	23.1
	+remove article	50.4	44.6	24.6

- higher error rates due to the rich morphology and free word order
- phrases are more important for this translation direction
- removing English articles is helpful

# Conclusions

- **an acceptable translation quality can be achieved with a very small amount of task-specific parallel text, especially if**
  - **conventional dictionaries and/or phrasal books**
  - **morpho-syntactic knowledge****are available**
  
- **translation systems built only on**
  - **conventional dictionary**
  - **phrasal book**
  - **extremely small parallel corpus****might be useful for document classification or multilingual information retrieval**

# References

- **Al-Onaizan & Germann<sup>+</sup>, 2000.**  
Translation with scarce resources.  
*17th National Conference on Artificial Intelligence (AAAI)*,  
pages 672–678, Austin, TX, August.
- **Callison-Burch & Osborne, 2003.**  
Co-training for statistical machine translation.  
*6th Annual CLUK Research Colloquium*, Edinburgh, UK, January.
- **Goldwater & McClosky, 2005**  
Improving statistical machine translation through morphological analysis.  
*Conf. on Empirical Methods for Natural Language Processing (EMNLP)*,  
Vancouver, Canada, October.
- **Lopez & Resnik, 2005.**  
Improved HMM alignment for languages with scarce resources.  
*ACL Workshop on Building and Using Parallel Texts: Data-Driven  
Machine Translation and Beyond*,  
pages 83–86, Ann Arbor, MI, June.

# References

- **Martin & Mihalcea<sup>+</sup>, 2005.**  
**Word alignments for languages with scarce resources.**  
***ACL Workshop on Building and Using Parallel Texts: Data-Driven Machine Translation and Beyond*,**  
**pages 65–74, Ann Arbor, MI, June.**
- **Matusov & Popović<sup>+</sup>, 2004.**  
**Statistical machine translation of spontaneous speech with scarce resources.**  
***Int. Workshop on Spoken Language Translation (IWSLT)*,**  
**pages 139–146, Kyoto, Japan, September.**
- **Niessen & Ney, 2004.**  
**Statistical machine translation with scarce resources using morpho-syntactic information.**  
***Computational Linguistics*, 30(2):181–204**

# References

- **Popović & Ney, 2005.**  
Exploiting phrasal lexica and additional morpho-syntactic language resources for statistical machine translation with scarce training data. *10th Annual Conf. on the European Association for Machine Translation (EAMT)*, pages 212–218, Budapest, Hungary, May.
- **Popović & Vilar<sup>+</sup>, 2005.**  
Augmenting a small parallel text with morpho-syntactic language resources for Serbian-English statistical machine translation. *ACL Workshop on Building and Using Parallel Texts: Data-Driven Machine Translation and Beyond*, pages 41–48, Ann Arbor, MI, June.
- **Popović & Ney, 2006.**  
POS-based word reorderings for statistical machine translation. *5th Int. Conf. on Language Resources and Evaluation (LREC)*, Genoa, Italy, May.