

**TC-Star OpenLab on Speech Translation  
Trento, 30/3 - 1/4/2006**

**Morpho-syntactic Information for Automatic Error Analysis of  
Statistical Machine Translation Output**

**Maja Popović**

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

# Overview

- **Introduction**
- **Morpho-syntactic Information and Automatic Evaluation**
- **Error Analysis**
- **Conclusion**

# Introduction

- standard automatic evaluation measures (WER, PER, BLEU, NIST) do not give any details about actual translation errors
  - subjective evaluation is time-consuming
- ⇒ use morpho-syntactic information in combination with automatic error measures

Popović & de Gispert<sup>+</sup> 06:

Morpho-syntactic Information for Automatic Error Analysis of  
Statistical Machine Translation Output

Submitted to *HLT/NAACL Workshop on Statistical Machine Translation 2006, FinTal 2006*

## Error Analysis

**Possible problems for translation of the Spanish-English language pair:**

- **reordering errors**

⇒ **compare WER and PER**  
**(large difference ↔ reordering errors)**

- **inflectional errors**

⇒ **compare PER of full forms and PER of base forms**  
**(large difference ↔ inflectional errors)**

# Corpus Statistics

## - EPPS Corpus -

		Spanish	English
<b>Training:</b> full	Sentences	<b>1281427</b>	
	Running Words+Punctuation	<b>36578514</b>	<b>34918192</b>
	Vocabulary	<b>153124</b>	<b>106496</b>
	Singletons [%]	<b>35.2</b>	<b>36.2</b>
<b>reduced</b>	Sentences	<b>13360</b>	
	Running Words+Punctuation	<b>385198</b>	<b>366055</b>
	Vocabulary	<b>22425</b>	<b>16326</b>
	Singletons [%]	<b>47.6</b>	<b>43.7</b>
<b>Develop:</b>	Sentences	<b>1008</b>	
	Running Words+Punctuation	<b>25778</b>	<b>26070</b>
	Distinct Words	<b>3895</b>	<b>3173</b>
	OOVs (full) [%]	<b>0.15</b>	<b>0.09</b>
	OOVs (reduced) [%]	<b>2.7</b>	<b>1.7</b>
<b>Test:</b>	Sentences	<b>840</b>	<b>1094</b>
	Running Words	<b>22774</b>	<b>26917</b>
	Distinct Words	<b>4081</b>	<b>3958</b>
	OOVs (full) [%]	<b>0.14</b>	<b>0.25</b>
	OOVs (running words) [%]	<b>2.8</b>	<b>2.6</b>

# Translation System

- **state-of-the-art translation system**
- **log-linear combination of seven models:**
  - phrase-based models (source to target and target to source)
  - single word based models at phrase level (source to target and target to source)
  - language model
  - phrase penalty and word penalty
- **results comparable to those obtained in the first TC-Star evaluation campaign**

# Error Analysis - Evaluation Results

## - reordering errors -

### relative difference between WER and PER

**English output:**

$1 - \frac{PER}{WER}$	full corpus		red. corpus	
	dev	test	dev	test
<b>nouns+adjectives</b>	<b>24.7</b>	<b>24.7</b>	<b>27.8</b>	<b>25.7</b>
<b>+reordering</b>	<b>21.6</b>	<b>20.8</b>	<b>21.2</b>	<b>20.1</b>
<b>verbs</b>	<b>4.9</b>	<b>4.1</b>	<b>4.9</b>	<b>4.6</b>
<b>adjectives</b>	<b>8.4</b>	<b>10.2</b>	<b>6.8</b>	<b>8.4</b>
<b>nouns</b>	<b>19.8</b>	<b>20.1</b>	<b>19.8</b>	<b>19.1</b>

**Spanish output:**

$1 - \frac{PER}{WER}$	full corpus		red. corpus	
	dev	test	dev	test
<b>nouns+adjectives</b>	<b>20.5</b>	<b>21.5</b>	<b>22.9</b>	<b>22.9</b>
<b>+reordering</b>	<b>18.9</b>	<b>20.3</b>	<b>20.6</b>	<b>19.8</b>
<b>verbs</b>	<b>3.2</b>	<b>3.3</b>	<b>3.4</b>	<b>3.9</b>
<b>adjectives</b>	<b>6.0</b>	<b>5.6</b>	<b>6.0</b>	<b>5.4</b>
<b>nouns</b>	<b>18.2</b>	<b>16.9</b>	<b>21.2</b>	<b>19.3</b>

# Error Analysis - Evaluation Results

## - inflectional errors -

### PER for different word classes

English output:

PER	full corpus		red. corpus	
	dev	test	dev	test
verbs	41.0	44.8	51.8	56.1
adjectives	28.2	27.3	38.2	38.1
nouns	22.6	23.0	39.2	31.7

Spanish output:

PER	full corpus		red. corpus	
	dev	test	dev	test
verbs	59.5	61.4	70.4	73.0
adjectives	40.4	41.8	50.0	50.9
nouns	27.8	28.5	35.0	37.0

# Error Analysis - Evaluation Results

## - inflectional errors -

### relative difference between base form PER and full form PER

Spanish output:

$1 - \frac{basePER}{fullPER}$	full corpus		red. corpus	
	dev	test	dev	test
verbs	25.9	26.9	25.9	23.7
adjectives	6.2	9.3	12.8	15.1
nouns	7.5	8.4	7.4	6.5

## Conclusion

- framework for automatic analysis of translation errors based on morpho-syntactic information
- results correspond to the results of the manual error analysis reported in [Vilar & Xu<sup>+</sup> 06]
- improvements of the baseline system adequately reflected on new measures

# Translation Results

## Spanish→English

<b>full corpus</b>	<b>dev</b>			<b>test</b>		
	<b>WER</b>	<b>PER</b>	<b>BLEU</b>	<b>WER</b>	<b>PER</b>	<b>BLEU</b>
<b>baseline</b>	<b>33.0</b>	<b>24.2</b>	<b>57.5</b>	<b>34.5</b>	<b>25.5</b>	<b>54.7</b>
<b>reorder adjective</b>	<b>32.4</b>	<b>23.9</b>	<b>58.3</b>	<b>33.5</b>	<b>25.2</b>	<b>56.4</b>

<b>reduced corpus</b>	<b>dev</b>			<b>test</b>		
	<b>WER</b>	<b>PER</b>	<b>BLEU</b>	<b>WER</b>	<b>PER</b>	<b>BLEU</b>
<b>baseline</b>	<b>39.2</b>	<b>28.4</b>	<b>48.7</b>	<b>41.8</b>	<b>30.7</b>	<b>43.2</b>
<b>reorder adjective</b>	<b>37.9</b>	<b>28.3</b>	<b>50.7</b>	<b>38.9</b>	<b>29.5</b>	<b>48.5</b>

# Translation Results

## English→Spanish

<b>full corpus</b>	<b>dev</b>			<b>test</b>		
	<b>WER</b>	<b>PER</b>	<b>BLEU</b>	<b>WER</b>	<b>PER</b>	<b>BLEU</b>
<b>baseline</b>	<b>39.8</b>	<b>30.2</b>	<b>50.5</b>	<b>39.7</b>	<b>30.6</b>	<b>47.8</b>
<b>reorder adjective</b>	<b>39.7</b>	<b>30.2</b>	<b>50.9</b>	<b>39.6</b>	<b>30.5</b>	<b>48.3</b>

<b>reduced corpus</b>	<b>dev</b>			<b>test</b>		
	<b>WER</b>	<b>PER</b>	<b>BLEU</b>	<b>WER</b>	<b>PER</b>	<b>BLEU</b>
<b>baseline</b>	<b>48.3</b>	<b>35.8</b>	<b>40.6</b>	<b>49.6</b>	<b>37.4</b>	<b>36.2</b>
<b>reorder adjective</b>	<b>47.4</b>	<b>35.6</b>	<b>41.7</b>	<b>48.1</b>	<b>36.5</b>	<b>37.7</b>

## References

- Banerjee & Lavie 05  
**METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments**  
In *Proc. ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan, June.
- Niessen & Ney 00  
Improving SMT quality with morpho-syntactic analysis.  
*18th International Conference on Computational Linguistics (CoLing)*  
pages 1081–1085, Saarbrücken, Germany, July.
- Niessen & Och<sup>+</sup> 00  
An Evaluation Tool for Machine Translation: Fast Evaluation for MT Research  
*Proc. 2nd Int. Conf. on Language Resources and Evaluation (LREC)*, pages 39–45, Athens, Greece, May.
- Niessen & Ney 04  
Statistical machine translation with scarce resources using morpho-syntactic information.  
*Computational Linguistics*, 30(2):181–204

## References

- Popović & Ney 06  
**POS-based Word Reorderings for Statistical Machine Translation**  
To appear: *5th Int. Conf. on Language Resources and Evaluation (LREC)*,  
Genoa, Italy, May.
- Popović & de Gispert<sup>+</sup> 06:  
**Morpho-syntactic Information for Automatic Error Analysis of**  
**Statistical Machine Translation Output**  
*Submitted to HLT/NAACL Workshop on Statistical Machine Translation 2006,*  
*FinTal 2006*
- Vilar & Xu<sup>+</sup> 06  
**Error Analysis of Statistical Machine Translation Output**  
To appear: *5th Int. Conf. on Language Resources and Evaluation (LREC)*,  
Genoa, Italy, May.