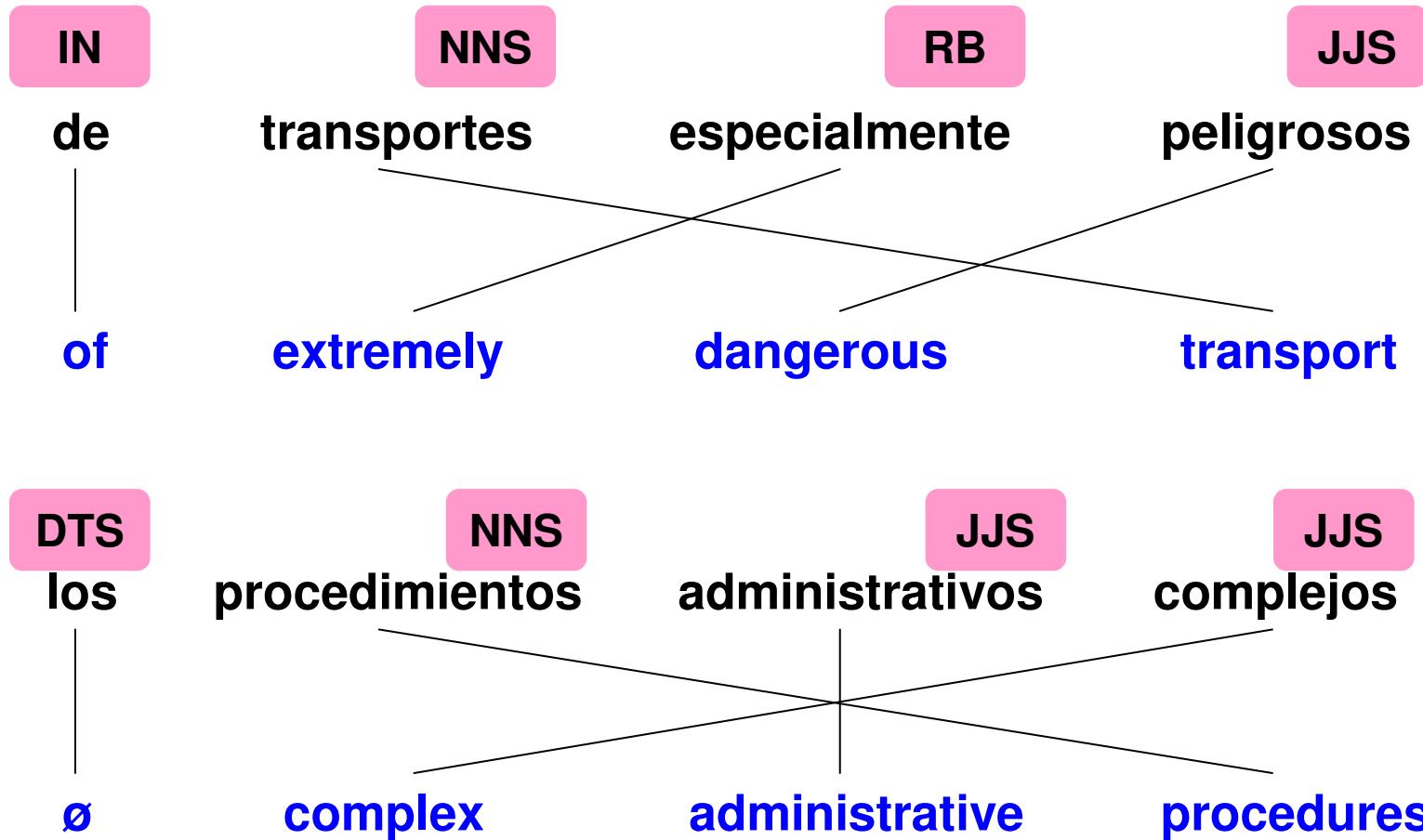


# **Morpho-Syntax in Statistical Machine Translation**

**Young-Suk Lee  
IBM T. J. Watson Research Center**

**OpenLab 2006  
March 30 – April 1, 2006**

# Reordering Rules: Motivations

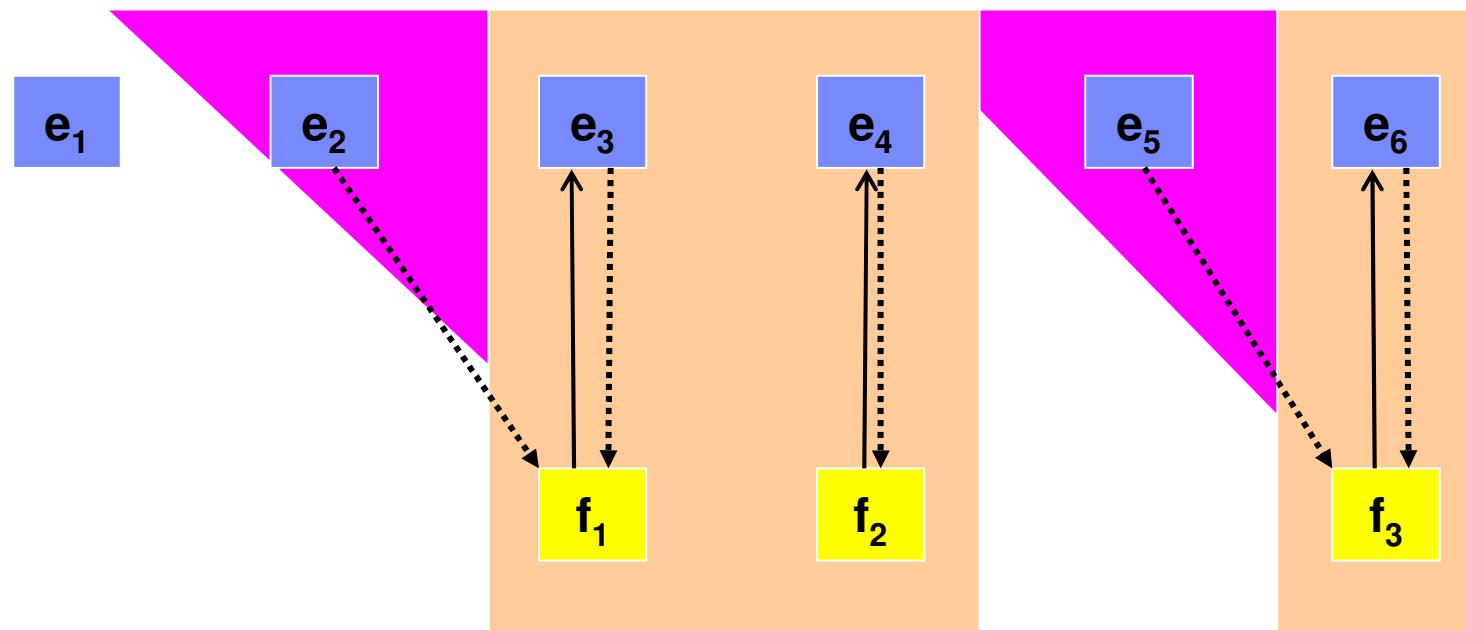


# Outline

- **Baseline Phrase Translation System**
  - Block Acquisition & Decoding
- **Acquisition of Reordering Rules**
  - Base Reordering Rules
  - Lexicalized Reordering Rules
- **Experimental Results**
- **Related and Ongoing Work**

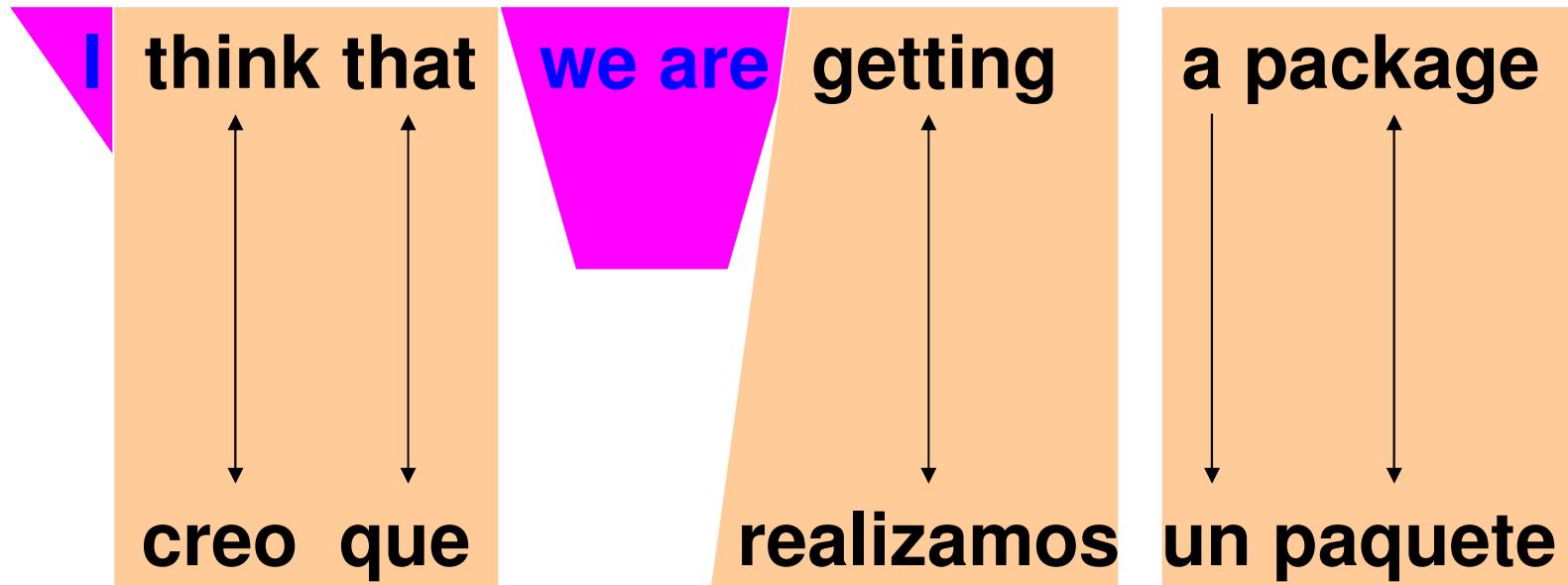
# Baseline Block Acquisition

**Block (b): a phrase translation pair consisting of source  $\bar{f}$  & target  $\bar{e}$  phrase**



Tillmann 2003, EMNLP Proceedings

# Extended Block Acquisition Algorithm



- o Expansion word list: A list of target words typically aligned to null source words (e.g. *I, we, are*)
- o Extend the target phrase to include an expansion word if it occurs in the neighborhood of a block

# Decoding

- **Phrase translation models**

- **Direct model:**

$$p(\bar{e} \mid \bar{f}) = \frac{\text{count}(\bar{e}, \bar{f})}{\sum_{\bar{e}'} \text{count}(\bar{e}', \bar{f})}$$

- **Source channel model:**

$$p(\bar{f} \mid \bar{e}) = \frac{\text{count}(\bar{f}, \bar{e})}{\sum_{\bar{f}'} \text{count}(\bar{f}', \bar{e})}$$

- **Block unigram model:**

$$p(b) = \frac{\text{count}(b)}{\sum_{b'} \text{count}(b')} , b = (\bar{e}, \bar{f})$$

## Decoding Cont'd ...

- IBM Model 1 cost per phrase in both directions

$$\sum_{j=1}^m -\log 10 \max_i p(f_j | e_i), 1 \leq i \leq n$$

- Word & part-of-speech tag trigram language models
- Word-level distortion models applied to blocks
  - Al-Onaizan 2004, DARPA MT Evaluation Workshop
- Word & block count penalty
  - Zens and Ney 2004, HLT Proceedings

# Acquisition of Base Reordering Rules

- Viterbi-align
  - Part-of-speech tagged source language corpus
  - Un-tagged target language corpus
- Identify the source language part-of-speech tag sequence (monotone increasing)
  - whose corresponding target word sequence is not monotone increasing
- Compute the reordering probabilities of each part-of-speech tag sequence

# Reordering Probability Computation

$$p(reorder_i | \overline{tag}_k) = \frac{\text{count}(reorder_i, \overline{tag}_k)}{\sum_{reorder'} \text{count}(reorder', \overline{tag}_k)}$$

DTS <sub>1</sub> NNS <sub>2</sub> JJS <sub>3</sub> JJS <sub>4</sub>	IN <sub>1</sub> NNS <sub>2</sub> RB <sub>3</sub> JJS <sub>4</sub>		
reorder'	p(reorder <sub>i</sub>   $\overline{tag}_k$ )	reorder'	p(reorder <sub>i</sub>   $\overline{tag}_k$ )
1 2 3 4	0.107	1 2 3 4	0.179
1 2 4 3	0.059	1 2 4 3	0.047
1 3 2 4	0.239	1 3 2 4	0.094
1 3 4 2	0.109	1 3 4 2	0.560
1 4 2 3	0.111	1 4 2 3	0.072
1 4 3 2	0.375	1 4 3 2	0.048

# One Best Reordering Rules

$$p(reorder_f | \overline{tag}) > p(reorder_s | \overline{tag}) + \alpha$$



<b>IN<sub>1</sub> NNS<sub>2</sub> RB<sub>3</sub> JJS<sub>4</sub></b>	<b>IN<sub>1</sub> RB<sub>3</sub> JJS<sub>4</sub> NNS<sub>2</sub></b>
<b>DTS<sub>1</sub> NNS<sub>2</sub> JJS<sub>3</sub> JJS<sub>4</sub></b>	<b>DTS<sub>1</sub> JJS<sub>4</sub> JJS<sub>3</sub> NNS<sub>2</sub></b>
<b>DT<sub>1</sub> NN<sub>2</sub> JJ<sub>3</sub> IN<sub>4</sub></b>	<b>DT<sub>1</sub> JJ<sub>3</sub> NN<sub>2</sub> IN<sub>4</sub></b>
<b>NNS<sub>1</sub> JJS<sub>2</sub> CC<sub>3</sub> JJS<sub>4</sub></b>	<b>JJS<sub>2</sub> CC<sub>3</sub> JJS<sub>4</sub> NNS<sub>1</sub></b>
<b>DT<sub>1</sub> NN<sub>2</sub> CC<sub>3</sub> NN<sub>4</sub> JJ<sub>5</sub></b>	<b>DT<sub>1</sub> JJ<sub>5</sub> NN<sub>2</sub> CC<sub>3</sub> NN<sub>4</sub></b>

# Lexicalization of Exceptions

el apoyo operativo de la<sub>1</sub>/DT Secretaría<sub>2</sub>/NN General<sub>3</sub>/JJ del<sub>4</sub>/IN Consejo

the operational support of the<sub>1</sub> Secretary<sub>2</sub> General<sub>3</sub> of<sub>4</sub> the Council

DT<sub>1</sub> NN<sub>2</sub> JJ<sub>3</sub>[~General] IN<sub>4</sub> → DT<sub>1</sub> JJ<sub>3</sub> NN<sub>2</sub> IN<sub>4</sub>

El Fondo, por supuesto, debe continuar cumpliendo con su misión de investigación sobre la búsqueda de<sub>1</sub>/IN variedades<sub>2</sub>/NNS más<sub>3</sub>/RB adaptadas<sub>4</sub>/JJS a la demanda y lo menos nocivas posible,

The Fund must of course continue to serve its purpose and pursue research into<sub>1</sub> varieties<sub>2</sub> more<sub>3</sub> suited<sub>4</sub> to demand and causing as little harm as possible .

IN<sub>1</sub> NNS<sub>2</sub> RB<sub>3</sub> JJS<sub>4</sub>[~adaptadas] → IN<sub>1</sub> RB<sub>3</sub> JJS<sub>3</sub> NNS<sub>2</sub>

# Lexicalized Reordering Rules

- Identify the key part-of-speech tag in the base reordering rules
- Replace the key part-of-speech tag with the corresponding word

o DT NN **JJ** IN → DT NN **General** IN

- Compute reordering probabilities of lexicalized part-of-speech tag sequences
- Exception word list
  - o If the reordering pattern with **the highest probability is monotone increasing**, select the word in the pattern as an exception

# Lexicalized Reordering Probabilities

DT <sub>1</sub> NN <sub>2</sub> General <sub>3</sub> IN <sub>4</sub>	<i>reorder</i> '	$p(reorder_i   tag_k)$
1 2 3 4	<b>reorder'</b>	<b>0.454</b>
1 2 4 3		<b>0.021</b>
1 3 2 4		<b>0.201</b>
1 3 4 2		<b>0.012</b>
1 4 2 3		<b>0.194</b>
1 4 3 2		<b>0.098</b>
4 1 2 3		<b>0.007</b>
4 1 3 2		<b>0.013</b>

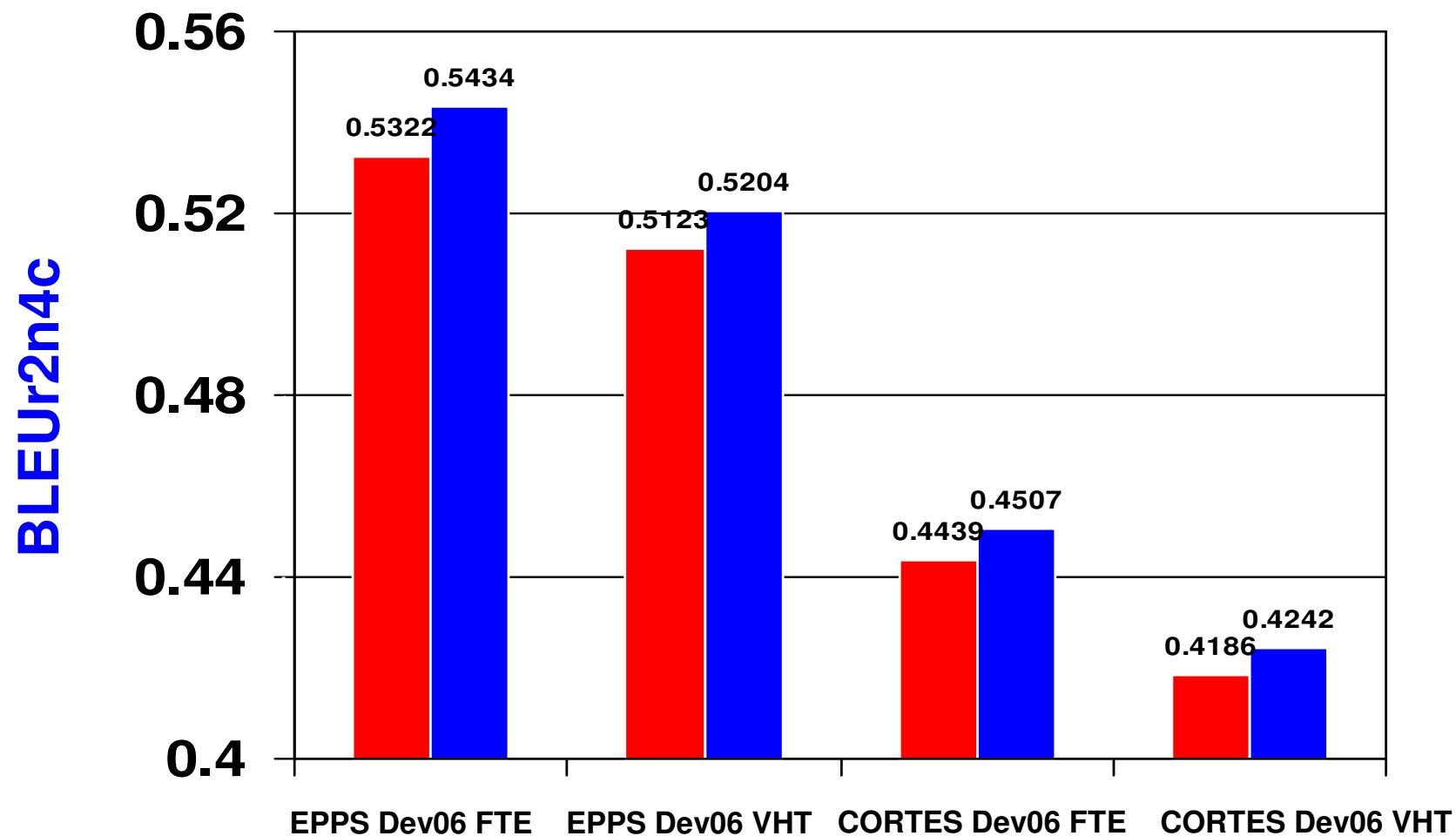
# Performance Evaluations

- Translation model training corpus
  - ~1.3 M sentence pairs from EPPS distributed by RWTH
- Language model training corpus
  - EPPS English corpus: ~35 M words
  - UN parallel corpus English (LDC94T4A): ~45 M words
  - English gigaword second edition (LDC2005T12): ~2.5 B words

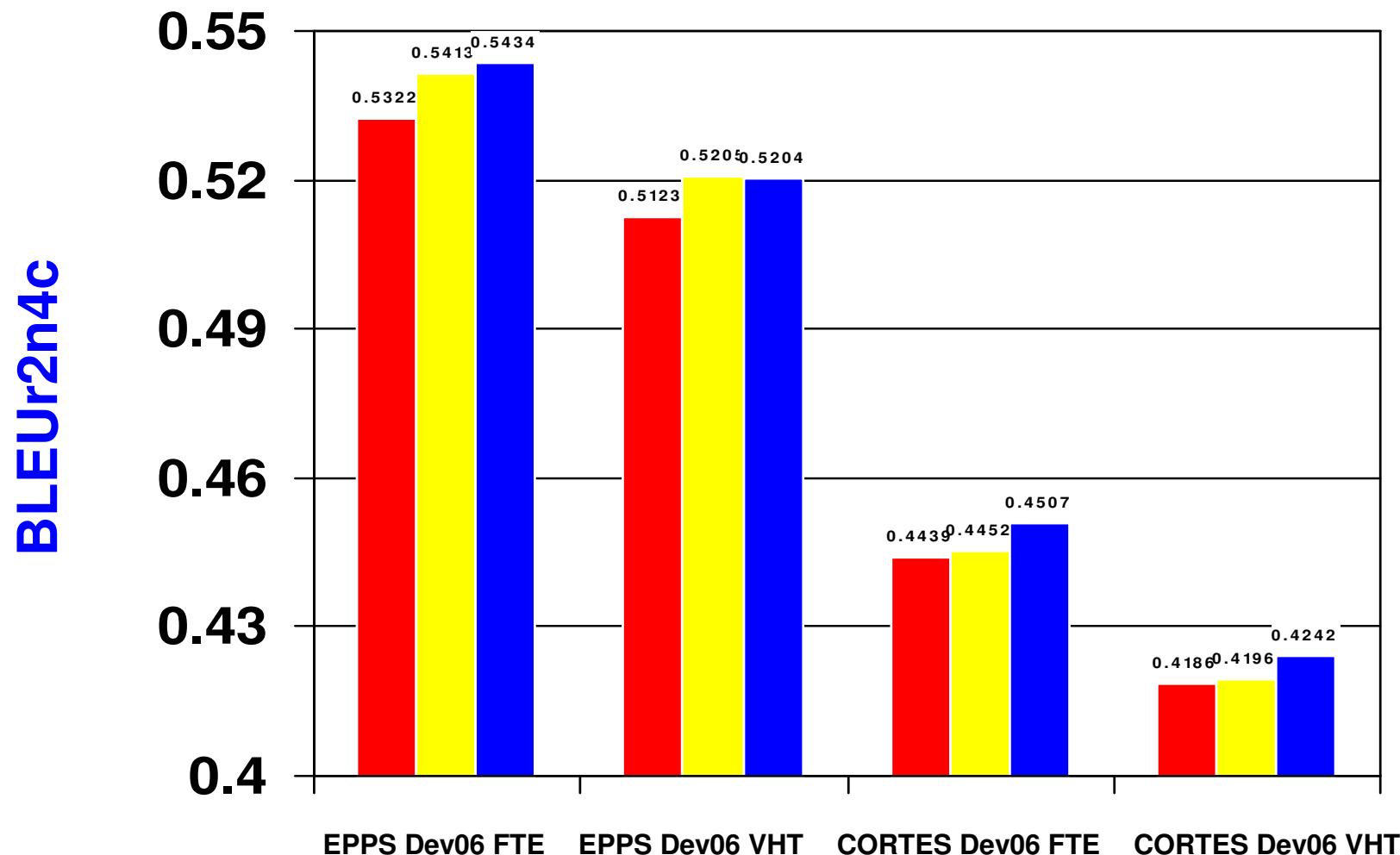
# Evaluation Corpus Statistics

Data Sets	# of Segments	Avg. Segment Length
EPPS Dev06 FTE	699	35 words/segment
EPPS Dev06 VHT	792	31 words/segment
CORTES Dev06 FTE	753	37 words/segment
CORTES Dev06 VHT	920	31 words/segment

# Lexicalized Reordering Rules: Impact



# Base vs. Lexicalized Reordering Rules



# Related Work

- N-best Reordering in Arabic-to-English Translation
  - Statistically significant performance improvement by applying local reordering to noun phrase parsed Arabic
  - *IBM Site Report: DARPA MT Evaluation Workshop 2004*
- Morphological Analysis for Statistical Machine Translation
  - Identify one to one word correspondences between Arabic and English to improve word to word translation qualities
  - *Companion Volume of HLT-NAACL 2004, pages 57–60*
- Local Reordering for Spanish-English Translations
  - Presentation at TC-STAR 2005 Evaluation Workshop
  - *April 21-22, 2005, Trento, Italy*

# Ongoing Work

- Non-local reordering models
  - [Se ha puesto a prueba]<sub>VP</sub> [su voluntad]<sub>NP</sub> →  
[Its will]<sub>NP</sub> [has been put to the test]<sub>VP</sub>
  - Todas sus Señorías firmaron [con los electores]<sub>PP</sub> [un contrato]<sub>NP</sub> → All your ladies and gentlemen signed [a contract]<sub>NP</sub> [with the electors]<sub>PP</sub>
- Integration of reordering models into the decoder