

**International Workshop on Spoken Language Translation  
Kyoto, Japan  
September 30 - October 1, 2004**

## **Alignment Templates: the RWTH SMT System**

**Oliver Bender, Richard Zens, Evgeny Matusov, and Hermann Ney**

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

# Content

- 1. overview: statistical machine translation**
- 2. loglinear models**
- 3. alignment templates**
- 4. feature functions**
- 5. minimum error training**
- 6.  $n$ -best lists and rescoring**
- 7. experimental results**
- 8. summary**

## Related work

- **F. J. Och and H. Ney. 2002. Discriminative training and maximum entropy models for statistical machine translation. In *Proc. of the 40th Annual Meeting of the Association for Computational Linguistics (ACL)*, pp. 295–302, Philadelphia, PA, July.**
- **F. J. Och, D. Gildea, S. Khudanpur, A. Sarkar, K. Yamada, A. Fraser, S. Kumar, L. Shen, D. Smith, K. Eng, V. Jain, Z. Jin, and D. Radev. 2004. A smorgasbord of features for statistical machine translation. In *Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics: HLT-NAACL 2004*, pp. 161–168, Boston, MA, May.**
- **A. Stolcke. 2002. SRILM - an extensible language modeling toolkit. In *Proc. Intl. Conf. Spoken Language Processing*, pp. 901–904, Denver, CO, September.**
- **D. Wu. 1997. Stochastic inversion transduction grammars and bilingual parsing of parallel corpora. *Computational Linguistics*, vol. 23, no. 3, pp. 377–403, September.**

# Overview: Statistical Machine Translation

- source string  $f_1^J = f_1 \dots f_j \dots f_J$  to be translated into a target string  $e_1^I = e_1 \dots e_i \dots e_I$ .
- classical source-channel approach:

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

- $Pr(f_1^J | e_1^I)$ : translation model  
(usually can be further decomposed into alignment and lexicon model)
- $Pr(e_1^I)$ : language model

## Loglinear models

- **alternative: direct modeling of the posterior probability  $Pr(e_1^I | f_1^J)$**
- **use a loglinear model (Och and Ney 2002):**

$$Pr(e_1^I | f_1^J) = p_{\lambda_1^M}(e_1^I | f_1^J) = \frac{\exp \left[ \sum_{m=1}^M \lambda_m h_m(e_1^I, f_1^J) \right]}{\sum_{e_1^I} \exp \left[ \sum_{m=1}^M \lambda_m h_m(e_1^I, f_1^J) \right]}$$

- **decision rule:**

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

- **advantages:**

- **easy integration of additional models/feature functions  $h_m$**
- **minimum error training of model scaling factors  $\lambda_m$**

# Alignment Templates

- **primary translation model: alignment templates**
- **describes the alignment between sequences of source and target words**
- **automatically trained word classes are used instead of words for better generalization**
- **translation model incorporates:**
  - **phrase alignment probability**
  - **probability to apply an alignment template**
  - **phrase translation probability**
- **alignment templates extracted automatically from automatic word alignments**



# Alignment Combination Heuristics

- **word alignments  $A_1$  and  $A_2$  are trained in source-to-target and target-to-source direction, respectively**
- **such alignments contain many-to-one mappings in one direction only**
- **alignment combination depends on the particular language pair**
- **best translation results achieved:**
  - **Chinese-English: using alignments which only allow many-to-one mappings of English words**
  - **Japanese-English: using “refined” alignments**
    - \* **extend intersection  $A_1 \cap A_2$  by additional points**
    - \* **add a new point if either a horizontal or a vertical direct neighbor point exists**



# Base Models Used in Search

- **alignment templates**
- **single-word translation model  $p(e|f)$**
- **word-based trigram language model**
- **class-based five-gram language model**
- **word penalty model**
- **phrase penalty model**
- **penalty for alignment template reorderings**

# Minimum Error Training

- optimize the model scaling factors  $\lambda_1^M$
- training criterion: minimal number of errors on a development corpus
- optimization with respect to a certain automatic translation score (100 – NIST, 1 – BLEU, WER)
- use the downhill simplex optimization algorithm
- translate the whole development corpus in each iteration of the algorithm
- algorithm converges after about 200 iterations

# Search

- **search characteristics:**
  - **reordering: within alignment templates: fixed in training**
  - **reordering of alignment templates: unconstrained or ITG (Japanese-English)**
  - **search organization along target string positions**
  - **beam search to handle the huge search space**
- **generation of  $n$ -best lists:**
  - **during search, generate word graphs**
  - **using the  $A^*$  search algorithm, compute  $n$ -best lists from the word graphs**

# Additional $n$ -best List Features

- **(inverse) IBM-1 lexicon model  $p(f|e)$  (as trained with GIZA++)**
  - + captures lexical co-occurrences, helpful for translation adequacy
- **deletion model**
  - + penalizes too short translation hypotheses
- **high-order  $n$ -gram language models ( $n = 4, 5, \dots, 9$ )**
  - + enrich the system with knowledge about longer target language phrases

# Deletion Model

- the produced translations are often shorter than the reference translations
- longer hypotheses are to be favored
- deletion model feature (Och et al. 2004): for a given threshold  $\alpha$ :
  - count the number of source words, for which the IBM-1 translation probability given any of the target words in the hypothesis is below  $\alpha$ .
  - use several features with different values of  $\alpha$  (0.1, 0.01, etc.)
- threshold  $\alpha$  tuned on a development corpus

## Experimental results

- IWSLT 2004 Evaluation
- rescoring improvements

# Evaluation Methodology

- **subjective evaluation as specified by the IWSLT 2004 consortium**
  - **translation fluency: from 1 (“incomprehensible”) to 5 (“flawless English”)**
  - **translation adequacy: how much information from a gold standard translation is contained in the hypothesis, from 1 (“none”) to 5 (“all”)**
- **objective evaluation: different automatic metrics computed using multiple references**
  - **Word Error Rate (mWER)**
  - **Position-Independent Word Error Rate (mPER)**
  - **BLEU score**
  - **NIST score**
  - **GTM score**

# BTEC Chinese-English Supplied Corpus Statistics

		Chinese	English
<b>train</b>	<b>sentences</b>	<b>20 000</b>	
	<b>words</b>	<b>182 904</b>	<b>160 523</b>
	<b>singletons</b>	<b>3 525</b>	<b>2 948</b>
	<b>vocabulary</b>	<b>7 643</b>	<b>6 982</b>
<b>dev</b>	<b>sentences</b>	<b>506</b>	
	<b>words</b>	<b>3 515</b>	<b>3 595</b>
<b>test</b>	<b>sentences</b>	<b>500</b>	
	<b>words</b>	<b>3 794</b>	<b>—</b>



# BTEC Japanese-English Supplied Corpus Statistics

		Japanese	English
<b>train</b>	<b>sentences</b>	<b>20 000</b>	
	<b>words</b>	<b>209 012</b>	<b>160 427</b>
	<b>singletons</b>	<b>4 108</b>	<b>2 956</b>
	<b>vocabulary</b>	<b>9 277</b>	<b>6 932</b>
<b>dev</b>	<b>sentences</b>	<b>506</b>	
	<b>words</b>	<b>4 374</b>	<b>3 595</b>
<b>test</b>	<b>sentences</b>	<b>500</b>	
	<b>words</b>	<b>4 370</b>	<b>—</b>

# BTEC Japanese-English Unrestricted Data Track Corpus Statistics

- additional resources:
  - full BTEC 1 Japanese-English corpus
  - Spoken Language Database (dialogs, hotel reservation domain)
- kindly provided by ATR

	Japanese	English
<b>train sentences</b>	<b>240 672</b>	
<b>words</b>	<b>1 974 407</b>	<b>1 770 190</b>
<b>singletons</b>	<b>8 975</b>	<b>3 658</b>
<b>vocabulary</b>	<b>26 037</b>	<b>14 301</b>
<b>dev sentences</b>	<b>506</b>	
<b>words</b>	<b>3 515</b>	<b>3 595</b>
<b>test sentences</b>	<b>500</b>	
<b>words</b>	<b>3 794</b>	<b>–</b>

# Official Evaluation Results

Language Pair	Data Track	Automatic Evaluation					Subj. Evaluation	
		mWER [%]	mPER [%]	BLEU [%]	NIST	GTM [%]	Fluency	Adequacy
<b>CE</b>	<b>Small</b>	<b>45.6</b>	<b>39.0</b>	<b>40.9</b>	<b>8.55</b>	<b>72.1</b>	<b>3.36</b>	<b>3.34</b>
<b>JE</b>	<b>Small</b>	<b>41.9</b>	<b>33.8</b>	<b>45.3</b>	<b>9.49</b>	<b>76.4</b>	<b>3.48</b>	<b>3.41</b>
	<b>Unrestricted</b>	<b>30.6</b>	<b>24.9</b>	<b>61.9</b>	<b>10.72</b>	<b>79.7</b>	<b>4.04</b>	<b>4.07</b>

- **balanced fluency/adequacy scores**
- **NIST score has the highest correlation with subjective ratings**

# Rescoring Improvements - Chinese-English

- error rates and scores on the development corpus (CSTAR 2003 test set)
- best overall performance achieved when optimizing the model scaling factors with respect to the NIST score
- base model scaling factors optimized using a narrow beam
- $n$ -best lists created using a broader beam
- each added feature results in performance gain

System	Error Rates		Accuracy Measures	
	mWER [%]	mPER [%]	BLEU [%]	NIST
baseline	55.2	45.6	34.8	7.76
broad beam	53.4	45.3	33.6	7.63
+ IBM-1 lexicon	50.9	42.1	36.4	8.06
+ deletion model	50.6	42.2	37.1	8.07
+ 9-gram LM	50.6	42.2	38.0	8.14

# Rescoring Improvements - Japanese-English

- error rates and scores on the development corpus (CSTAR 2003 test set)
- ITG reordering constraints in search improve the translation quality

System	Error Rates		Accuracy Measures	
	mWER [%]	mPER [%]	BLEU [%]	NIST
<b>baseline</b>	<b>48.7</b>	<b>38.6</b>	<b>44.3</b>	<b>9.10</b>
<b>+ ITG constraints</b>	<b>45.1</b>	<b>36.0</b>	<b>47.3</b>	<b>9.32</b>
<b>+ broad beam</b>	<b>49.5</b>	<b>37.3</b>	<b>45.0</b>	<b>9.32</b>
<b>+ IBM-1 lexicon</b>	<b>44.6</b>	<b>35.7</b>	<b>48.9</b>	<b>9.71</b>
<b>+ deletion model</b>	<b>43.2</b>	<b>34.7</b>	<b>50.1</b>	<b>9.80</b>
<b>+ 5-gram LM</b>	<b>42.6</b>	<b>34.2</b>	<b>51.5</b>	<b>9.92</b>

# Conclusions

- translation system based on loglinear model combination
- additional knowledge sources easily integrated as features
- phrasal context and local word reorderings are important
  - ⇒ captured in the alignment templates model
- direct optimization of base models using minimum error training of model scaling factors
- an additional deletion model feature penalizes too short translations
- scaling factors for additional features optimized using  $n$ -best lists of translation hypotheses
- optimization of the RWTH system with respect to the NIST score seems to correspond best to subjective evaluation criteria
- on the BTEC Chinese-English and Japanese-English tasks, translations of good quality were produced