

# Rich Source-Side Context for Statistical Machine Translation

Kevin Gimpel and Noah A. Smith

Language Technologies Institute  
School of Computer Science  
Carnegie Mellon University  
Pittsburgh, PA 15213, USA  
{kgimpel, nasmith}@cs.cmu.edu

## Abstract

We explore the augmentation of statistical machine translation models with features of the *context* of each phrase to be translated. This work extends several existing threads of research in statistical MT, including the use of context in example-based machine translation (Carl and Way, 2003) and the incorporation of word sense disambiguation into a translation model (Chan et al., 2007). The context features we consider use surrounding words and part-of-speech tags, local syntactic structure, and other properties of the source language sentence to help predict each phrase’s translation. Our approach requires very little computation beyond the standard phrase extraction algorithm and scales well to large data scenarios. We report significant improvements in automatic evaluation scores for Chinese-to-English and English-to-German translation, and also describe our entry in the WMT-08 shared task based on this approach.

## 1 Introduction

Machine translation (MT) by statistical modeling of bilingual phrases is one of the most successful approaches in the past few years. Phrase-based MT systems are straightforward to train from parallel corpora (Koehn et al., 2003) and, like the original IBM models (Brown et al., 1990), benefit from standard language models built on large monolingual, target-language corpora (Brants et al., 2007). Many of these systems perform well in competitive evaluations and scale well to large-data situations

(NIST, 2006; Callison-Burch et al., 2007). End-to-end phrase-based MT systems can be built entirely from freely-available tools (Koehn et al., 2007).

We follow the approach of Koehn et al. (2003), in which we translate a source-language sentence  $f$  into the target-language sentence  $\hat{e}$  that maximizes a linear combination of features and weights:<sup>1</sup>

$$\langle \hat{e}, \hat{\mathbf{a}} \rangle = \operatorname{argmax}_{\langle \mathbf{e}, \mathbf{a} \rangle} \operatorname{score}(\mathbf{e}, \mathbf{a}, \mathbf{f}) \quad (1)$$

$$= \operatorname{argmax}_{\langle \mathbf{e}, \mathbf{a} \rangle} \sum_{m=1}^M \lambda_m h_m(\mathbf{e}, \mathbf{a}, \mathbf{f}) \quad (2)$$

where  $\mathbf{a}$  represents the segmentation of  $e$  and  $f$  into phrases and a correspondence between phrases, and each  $h_m$  is a  $\mathbb{R}$ -valued feature with learned weight  $\lambda_m$ . The translation is typically found using beam search (Koehn et al., 2003). The weights  $\langle \lambda_1, \dots, \lambda_M \rangle$  are typically learned to directly minimize a standard evaluation criterion on development data (e.g., the BLEU score; Papineni et al., (2002)) using numerical search (Och, 2003).

Many features are used in phrase-based MT, but nearly ubiquitous are estimates of the conditional translation probabilities  $p(e_i^j | f_k^l)$  and  $p(f_k^l | e_i^j)$  for each phrase pair  $\langle e_i^j, f_k^l \rangle$  in the candidate sentence pair.<sup>2</sup> In this paper, we add and evaluate fea-

<sup>1</sup>In the statistical MT literature, this is often referred to as a “log-linear model,” but since the score is normalized during neither parameter training nor decoding, and is never interpreted as a log-probability, it is essentially a linear combination of feature functions. Since many of the features are actually probabilities, this linear combination is closer to a *mixture model*.

<sup>2</sup>We will use  $x_i^j$  to denote the subsequence of  $x$  containing the  $i$ th through  $j$ th elements of  $x$ , inclusive.

tures that condition on additional context features on the *source* ( $f$ ) side:

$$p(e_i^j \mid \text{Phrase} = f_k^\ell, \text{Context} = \langle f_1^{k-1}, f_{\ell+1}^F, \dots \rangle)$$

The advantage of considering context is well-known and exploited in the example-based MT community (Carl and Way, 2003). Recently researchers have begun to use source phrase context information in statistical MT systems (Stroppa et al., 2007). Statistical NLP researchers understand that conditioning a probability model on more information is helpful only if there are sufficient training data to accurately *estimate* the context probabilities.<sup>3</sup> Sparse data are often the death of elaborate models, though this can be remedied through careful smoothing.

In this paper we leverage the existing linear model (Equation 2) to bring source-side context into phrase-based MT in a way that is robust to data sparseness. We interpret the linear model as a mixture of many probability estimates based on different context features, some of which may be very sparse. The mixture coefficients are trained in the usual way (“minimum error-rate training,” Och, 2003), so that the additional context is exploited when it is useful and ignored when it isn’t.

The paper proceeds as follows. We first review related work that enriches statistical translation models using context (§2). We then propose a set of source-side features to be incorporated into the translation model, including the novel use of syntactic context from source-side parse trees and global position within  $f$  (§3). We explain why analogous *target*-side features pose a computational challenge (§4). Specific modifications to the standard training and evaluation paradigm are presented in §5. Experimental results are reported in §6.

## 2 Related Work

Stroppa et al. (2007) added source-side context features to a phrase-based translation system, including conditional probabilities of the same form that we use. They consider up to two words and/or POS tags of context on either side. Because of the aforementioned data sparseness problem, they use a decision-

<sup>3</sup>An illustrative example is the debate over the use of bilingualized grammar rules in statistical parsing (Gildea, 2001; Bikel, 2004).

tree classifier that implicitly smooths relative frequency estimates. The method improved over a standard phrase-based baseline trained on small amounts of data (< 50K sentence pairs) for Italian → English and Chinese → English. We explore a significantly larger space of context features, a smoothing method that more naturally fits into the widely used, error-driven linear model, and report a more comprehensive experimental evaluation (including feature comparison and scaling up to very large datasets).

Recent research on the use of word-sense disambiguation in machine translation also points toward our approach. For example, Vickrey et al. (2005) built classifiers inspired by those used in word sense disambiguation to fill in blanks in a partially-completed translation. Giménez and Márquez (2007) extended the work by considering phrases and moved to full translation instead of filling in target-side blanks. They trained an SVM for each source language phrase using local features of the sentences in which the phrases appear. Carpuat and Wu (2007) and Chan et al. (2007) embedded state-of-the-art word sense disambiguation modules into statistical MT systems, achieving performance improvements under several automatic measures for Chinese → English translation.

Our approach is also reminiscent of example-based machine translation (Nagao, 1984; Somers, 1999; Carl and Way, 2003), which has for many years emphasized use of the context in which source phrases appear when translating them. Indeed, like the example-based community, we do not begin with any set of assumptions about *which kinds* of phrases require additional disambiguation (cf. the application of word-sense disambiguation, which is motivated by lexical ambiguity). Our feature-rich approach is omnivorous and can exploit *any* linguistic analysis of an input sentence.

## 3 Source-Side Context Features

Adding features to the linear model (Equation 2) that consider more of the source sentence requires changing the decoder very little, if at all. The reason is that the source sentence is fully observed, so the information to be predicted is the same as before—the difference is that we are using more clues to carry out the prediction.

We see this as an opportunity to include many more features in phrase-based MT without increasing the cost of decoding at runtime. This discussion is reminiscent of an advantage gained by moving from hidden Markov models to conditional random fields for sequence labeling tasks. While the same core algorithm is used for decoding with both models, a CRF allows inclusion of features that consider the *entire* observed sequence—i.e., more of the observable context of each label to be predicted. Although this same advantage was already obtained in statistical MT through the transition from “noisy channel” translation models to (log-)linear models, the customary set of features used in most phrase-based systems does not take full advantage of the observed data.

The standard approach to estimating the phrase translation conditional probability features is via relative frequencies (here  $e$  and  $f$  are phrases):

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

Our new features all take the form  $p(e | f, f_{\text{context}})$ , where  $e$  is the target language phrase,  $f$  is the source language phrase, and  $f_{\text{context}}$  is the context of the source language phrase in the sentence in which it was observed. Like the context-bare conditional probabilities, we estimate probability features using relative frequencies:

$$p(e | f, f_{\text{context}}) = \frac{\text{count}(e, f, f_{\text{context}})}{\sum_{e'} \text{count}(e', f, f_{\text{context}})}$$

Since we expect that adding conditioning variables will lead to sparser counts and therefore more zero estimates, we compute features for many different types of context. To combine the many differently-conditioned features into a single model, we provide them as features to the linear model (Equation 2) and use minimum error-rate training (Och, 2003) to obtain interpolation weights  $\lambda_m$ . This is similar to an interpolation of backed-off estimates, if we imagine that all of the different contexts are differently-backed off estimates of the *complete* context. The error-driven weight training effectively smooths one implicit context-rich estimate  $p(e | f, f_{\text{context}})$  so that all of the backed-off es-

timates are taken into account, including the original  $p(e | f)$ . Our approach is asymmetrical; we have not, for example, estimated features of the form  $p(f, f_{\text{context}} | e)$ .

We next discuss the specific source-side context features used in our model.

### 3.1 Lexical Context Features

The most obvious kind of context of a source phrase  $f_k^\ell$  is the  $m$ -length sequence before it ( $f_{k-m}^{k-1}$ ) and the  $m$ -length sequence after it ( $f_{\ell+1}^{\ell+m}$ ). We include context features for  $m \in \{1, 2\}$ , padding sentences with  $m$  special symbols at the beginning and at the end. For each value of  $m$ , we include three features:

- $p(e | f, f_{k-m}^{k-1})$ , the left lexical context;
- $p(e | f, f_{\ell+1}^{\ell+m})$ , the right lexical context;
- $p(e | f, f_{k-m}^{k-1}, f_{\ell+1}^{\ell+m})$ , both sides.

### 3.2 Shallow Syntactic Features

Lexical context features, especially when  $m > 1$ , are expected to be sparse. Representing the context by part-of-speech (POS) tags is one way to overcome that sparseness. We used the same set of the lexical context features described above, but with POS tags replacing words in the context. We also include a feature which conditions on the POS tag sequence of the actual phrase being translated.

### 3.3 Syntactic Features

If a robust parser is available for the source language, we can include context features from parse trees. We used the following parse tree features:

- Is the phrase (exactly) a constituent?
- What is the nonterminal label of the lowest node in the parse tree that covers the phrase?
- What is the nonterminal label or POS of the highest nonterminal node that ends immediately before the phrase? Begins immediately after the phrase?
- Is the phrase strictly to the left of the root word, does it contain the root word, or is it strictly to the right of the root word? (Requires a parse with head annotations.)

We also used a feature that conditions on both features in the third bullet point above.

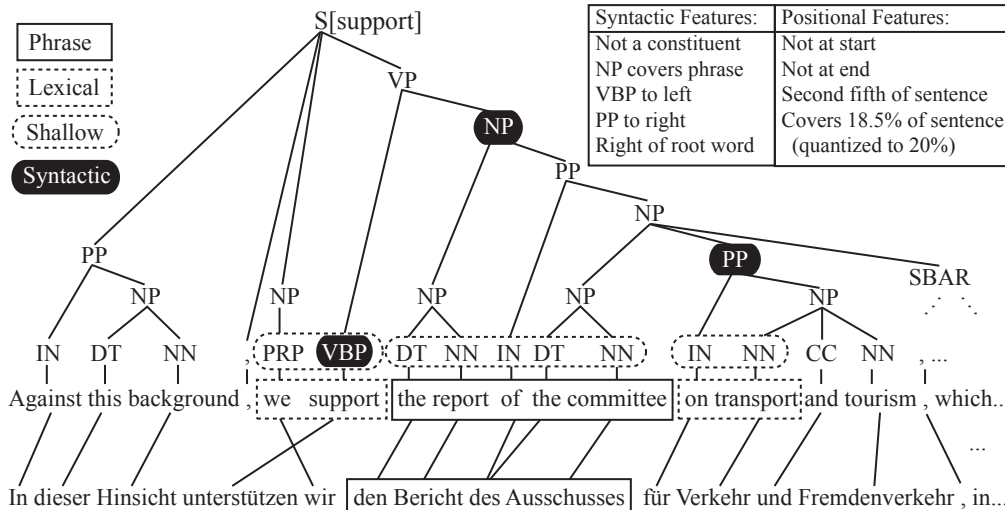


Figure 1: A (partial) sentence pair from the WMT-07 Europarl training corpus. Processing of the data (parsing, word alignment) was done as discussed in §6. The phrase pair of interest is boxed and context features are shown in dotted shapes. The context features help determine whether the phrase should be translated as “der Bericht des Ausschusses” (nominative case) or “den Bericht des Ausschusses” (accusative case). See text for details.

### 3.4 Positional Features

We include features based on the position of the phrase in the source sentence, the phrase length, and the sentence length. These features use information from the entire source sentence, but are not syntactic. For a phrase  $f_k^\ell$  in a sentence  $f$  of length  $n$ :

- Is the phrase at the start of the sentence ( $k = 1$ )?
- Is the phrase at the end of the sentence ( $\ell = n$ )?
- A quantization of  $r = \frac{k + \frac{\ell - k + 1}{2}}{n}$ , the relative position in  $(0, 1)$  of the phrase’s midpoint within  $f$ . We choose the smallest  $q \in \{0.2, 0.4, 0.6, 0.8, 1\}$  such that  $q > r$ .
- A quantization of  $c = \frac{\ell - k + 1}{n}$ , the fraction of the words in  $f$  that are covered by the phrase. We choose the smallest  $q \in \{\frac{1}{40}, \frac{1}{20}, \frac{1}{10}, \frac{1}{5}, \frac{1}{3}, 1\}$  such that  $q > c$ .

An illustration of the context features is shown in Fig. 1. Consider the phrase pair “the report of the committee”/“den Bericht des Ausschusses” extracted by our English  $\rightarrow$  German baseline MT system (described in §6.3). The German word “Bericht” is a masculine noun; therefore, it takes the article “der” in the nominative case, “den” in the accusative case, and “dem” in the dative case. These three translations are indeed available in the phrase table for “the report of the committee” (see Table 1, “no context” column), with relatively high entropy.

The choice between “den” and “der” must be made by the language model alone.

Knowing that the phrase follows a verb, or appears to the right of the sentence’s root word, or within the second fifth of the sentence should help. Indeed, a probability distribution that conditions on context features gives more peaked distributions that give higher probability to the correct translation, given *this* context, and lower probability given some *other* contexts (see Table 1).

### 4 Why Not Target-Side Context?

While source context is straightforward to exploit in a model, including target-side context features breaks one of the key independence assumptions made by phrase-based translation models: the translations of the source-side phrases are conditionally *independent* of each other, given  $f$ , thereby requiring new algorithms for decoding (Marino et al., 2006).

We suggest that target-side context may already be well accounted for in current MT systems. Indeed, *language models* pay attention to the local context of phrases, as do reordering models. The recent emphasis on improving these components of a translation system (Brants et al., 2007) is likely due in part to the widespread availability of NLP tools for the language that is most frequently the target: English. We will demonstrate that NLP tools (tag-

$g$	no context	Shallow: 2 POS on left		Syntax: _ of root		Positional: rel. pos.	
		*"PRP VBP"	"VBN IN"	*right	left	*2nd fifth	1st fifth
den bericht des ausschusses	0.3125	<b>1.0000</b>	0.3333	<b>0.5000</b>	0.0000	<b>0.6000</b>	0.0000
der bericht des ausschusses	0.3125	0.0000	0.0000	0.1000	<b>0.6667</b>	0.2000	<b>0.6667</b>
dem bericht des ausschusses	0.2500	0.0000	<b>0.6667</b>	0.3000	0.1667	0.0000	0.1667

Table 1: Phrase table entries for “the report of the committee” and their scores under different contexts. These are the top three phrases in the baseline English  $\rightarrow$  German system (“no context” column). Contexts from the source sentence in Fig. 1 (starred) predict correctly; we show also alternative contexts that give very different distributions.

gers and parsers) for the *source* side can be used to improve the translation model, exploiting analysis tools for other languages.

## 5 Implementation

The additional data required to compute the context features is extracted along with the phrase pairs during execution of the standard phrase extraction algorithm, affecting phrase extraction and scoring time by a constant factor. We avoid the need to modify the standard phrase-based decoder to handle context features by appending a unique identifier to each token in the sentences to be translated. Then, we precompute a phrase table for the phrases in these sentences according to the phrase contexts. To avoid extremely long lists of translations of common tokens, we filter the generated phrase tables, removing entries for which the estimate of  $p(e | f) < c$ , for some small  $c$ . In our experiments, we fixed  $c = 0.0002$ . This filtering reduced time for experimentation dramatically and had no apparent effect on the translation output. We did not perform any filtering for the baseline system.

## 6 Experiments

In this section we present experimental results using our context-endowed phrase translation model with a variety of different context features, on Chinese  $\rightarrow$  English, German  $\rightarrow$  English, and English  $\rightarrow$  Ger-

Context features	Chinese $\rightarrow$ English (UN)		
	BLEU	NIST	METEOR
None	0.3715	7.918	0.6486
Lexical	<b>0.4030</b>	<b>8.367</b>	<b>0.6716</b>
Shallow	<b>0.3807</b>	<b>7.981</b>	<b>0.6523</b>
Lexical + Shallow	<b>0.4030</b>	<b>8.403</b>	<b>0.6703</b>
Syntactic	<b>0.3823</b>	<b>7.992</b>	<b>0.6531</b>
Positional	<b>0.3775</b>	7.958	<b>0.6510</b>

Table 2: Chinese  $\rightarrow$  English experiments: training and testing on UN data. Boldface marks scores significantly higher than “None.”

man translation tasks. Dataset details are given in Appendices A (Chinese) and B (German).

**Baseline** We use the Moses MT system (Koehn et al., 2007) as a baseline and closely follow the example training procedure given for the WMT-07 and WMT-08 shared tasks.<sup>4</sup> In particular, we perform word alignment in each direction using GIZA++ (Och and Ney, 2003), apply the “grow-diag-final-and” heuristic for symmetrization and use a maximum phrase length of 7. In addition to the two phrase translation conditionals  $p(e | f)$  and  $p(f | e)$ , we use lexical translation probabilities in each direction, a word penalty, a phrase penalty, a length-based reordering model, a lexicalized reordering model, and an  $n$ -gram language model, SRILM implementation (Stolcke, 2002) with modified Kneser-Ney smoothing (Chen and Goodman, 1998). Minimum error-rate (MER) training (Och, 2003) was applied to obtain weights ( $\lambda_m$  in Equation 2) for these features. A recaser is trained on the target side of the parallel corpus using the script provided with Moses. All output is recased and detokenized prior to evaluation.

**Evaluation** We evaluate translation output using three automatic evaluation measures: BLEU (Papineni et al., 2002), NIST (Doddington, 2002), and METEOR (Banerjee and Lavie, 2005, version 0.6).<sup>5</sup> All measures used were the case-sensitive, corpus-level versions. The version of BLEU used was that provided by NIST. Significance was tested using a paired bootstrap (Koehn, 2004) with 1000 samples ( $p < 0.05$ ).<sup>6</sup>

<sup>4</sup><http://www.statmt.org/wmt08>

<sup>5</sup>METEOR details: For English, we use exact matching, Porter stemming, and WordNet synonym matching. For German, we use exact matching and Porter stemming. These are the same settings that were used to evaluate systems for the WMT-07 shared task.

<sup>6</sup>Code implementing this test for these metrics can be freely downloaded at <http://www.ark.cs.cmu.edu/MT>.

Context features	Chinese → English					
	Testing on UN			Testing on News (NIST 2005)		
	BLEU	NIST	METEOR	BLEU	NIST	METEOR
Training on in-domain data only:						
None	0.3715	7.918	0.6486	0.2700	7.986	0.5314
Training on all data:						
None	0.3615	7.797	0.6414	0.2593	7.697	0.5200
Lexical	<b>0.3898</b>	<b>8.231</b>	<b>0.6697</b>	0.2522	7.852	0.5273
Shallow: $\leq 1$ POS tag	0.3611	7.713	0.6430	0.2669	<b>8.243</b>	<b>0.5526</b>
Shallow: $\leq 2$ POS tags	0.3657	7.808	0.6455	0.2591	7.843	0.5288
Lexical + Shallow	<b>0.3886</b>	<b>8.245</b>	<b>0.6675</b>	0.2628	7.881	0.5290
Syntactic	0.3717	7.899	<b>0.6531</b>	0.2653	<b>8.123</b>	<b>0.5403</b>
Lexical + Syntactic	<b>0.3926</b>	<b>8.224</b>	<b>0.6636</b>	0.2572	7.774	0.5234
Positional	0.3647	7.766	0.6469	0.2648	7.891	0.5275
All	<b>0.3772</b>	<b>8.176</b>	<b>0.6582</b>	0.2566	7.775	0.5225
Feature selection (see Sec. 6.4)	<b>0.3843</b>	<b>8.079</b>	<b>0.6594</b>	0.2730	8.059	0.5343

Table 3: Chinese → English experiments: first row shows baseline performance when training only on in-domain data for each task; all other rows show results when training on *all* data (UN and News). Left half shows results when tuning and testing on UN test sets; right half shows results when tuning on NIST 2004 News test set and testing on NIST 2005. For feature selection, an additional set of unseen data was used: 2000 held-out sentences from the UN data for the left half and the NIST 2003 test set for the right half. Boldface marks scores that are significantly higher than the first row, in-domain baseline.

## 6.1 Chinese → English

For our Chinese → English experiments, two kinds of data were used: UN proceedings, and newswire as used in NIST evaluations.

**UN Data** UN data results are reported in Table 2. Significant improvements are obtained on all three evaluation measures—e.g., more than 3 BLEU points—using lexical or lexical and shallow features. While improvements are smaller for other features and feature combinations, performance is not *harmed* by conditioning on context features. Note that using syntactic features gave 1 BLEU point of improvement.

**News Data** In News data experiments, none of our features obtained BLEU performance statistically distinguishable from the baseline of 0.2700 BLEU (neither better, nor worse). The News training corpus is less than half the size of the UN training corpus (in words); unsurprisingly, the context features were too sparse to be helpful. Further, newswire are less formulaic and repetitive than UN proceedings, so contexts do not generalize as well from training to test data. Fortunately, our “error-minimizing mixture” approach protects the BLEU score, which the  $\lambda_m$  are tuned to optimize.

**Combined UN + News Data** Our next experiment used *all* of the available training data ( $> 200M$  words on each side) to train the models, in-domain  $\lambda_m$  tuning, and testing for each domain separately; see Table 3. Without context features, training on mixed-domain data consistently *harms* performance. With contexts that include lexical features, the mixed-domain model significantly outperforms the *in-domain* baseline for UN data. These results suggest that context features enable better use of out-of-domain data, an important advantage for statistical MT since parallel data often arise from very different sources than those of “real-world” translation scenarios. On News data, context features did not give a significant advantage on the BLEU score, though syntactic and  $\leq 1$  POS contexts did give significant NIST and METEOR improvements over the in-domain baseline.

## 6.2 German → English

We do not report full results for this task, because the context features neither helped nor hurt performance significantly. We believe this is due to data sparseness resulting from the size of the training corpus (26M German words), German’s relatively rich morphology, and the challenges of German parsing.

Context features	English → German		
	BLEU	NIST	METEOR
None	0.2069	6.020	0.2811
Lexical	0.2018	6.031	0.2772
Shallow	0.2017	5.911	0.2748
Syntactic	0.2077	6.049	0.2829
Positional	0.2045	5.930	0.2772
Lex. + Shal. + Syn.	0.2045	<b>6.061</b>	0.2817
All	0.2053	6.009	0.2797
Feature selection	0.2080	6.009	0.2807

Table 4: English → German experiments: training and testing on Europarl data. WMT-07 Europarl parallel training data was used for training, dev06 was used for tuning, devtest06 was used for feature selection and developmental testing, and test07 was used for final testing. Boldface marks scores significantly higher than “None.”

### 6.3 English → German

English → German results are shown in Table 4. The baseline system here is highly competitive, having scored higher on automatic evaluation measures than any other system in the WMT-07 shared task (Callison-Burch et al., 2007). Though most results are not statistically significant, small improvements do tend to come from *syntactic* context features. Comparing with the German → English experiment, we attribute this effect to the high accuracy of the English parser compared to the German parser.

### 6.4 Feature Selection

Translation performance does not always increase when features are added to the model. This motivates the use of feature selection algorithms to choose a subset of features to optimize performance. We experimented with several feature selection algorithms based on information-theoretic quantities computed among the source phrase, the target phrase, and the context, but found that a simple forward variable selection algorithm (Guyon and Elisseeff, 2003) worked best. In this procedure, we start with no context features and, at each iteration, add the single feature that results in the largest increase in BLEU score on an unseen development set after  $\lambda_m$  tuning. The algorithm terminates if no features are left or if none result in an increase in BLEU. We ran this algorithm to completion for the two Chinese → English tune/test sets (training on *all* data in each case) and the English → German task;

see Tables 3 and 4. In all cases, the algorithm finishes after  $\leq 4$  iterations.

Feature selection for Chinese → English (UN) first chose the lexical feature “1 word on each side,” then the positional feature indicating which fifth of the sentence contains the phrase, and finally the lexical feature “1 word on right.” For News, the features chosen were the shallow syntactic feature “1 POS on each side,” then the positional beginning-of-sentence feature, then the position relative to the root (a syntactic feature). For English → German, the shallow syntactic feature “2 POS on left,” then the lexical feature “1 word on right” were selected.

In the case where context features were most helpful (Chinese → English UN data), we found feature selection to be competitive at 2.28 BLEU points above the no-context baseline, but not the best achieved. In the other two cases (Chinese → English News and English → German Europarl), our best results were achieved using these automatically selected features, and in the Chinese → English News case, improvements on all three scores (including 1.37 BLEU points) are significant compared to the no-context baseline trained on the same data.

### 6.5 WMT-08 Shared Task: English → German

Since we began this research before the release of the data for the WMT-08 shared task, we performed the majority of our experiments using the data released for the WMT-07 shared task (see Appendix B). To prepare our entry for the 2008 shared task, we trained a baseline system on the 2008 data using a nearly identical configuration.<sup>7</sup> Table 5 compares performance of the baseline system (with no context features) to performance with the two context features chosen automatically as described in §6.4. In addition to the devtest06 data, we report results on the 2007 and 2008 Europarl test sets. Most improvements were statistically significant.

## 7 Future Work

In future work, we plan to apply more sophisticated learning algorithms to rich-feature phrase table estimation. Context features can also be used as conditioning variables in other components of translation

<sup>7</sup>The only differences were the use of a larger max sentence length threshold of 55 tokens instead of 50, and the use of the better-performing “englishFactored” Stanford parser model.

System	devtest06			test07			test08		
	BLEU	NIST	METEOR	BLEU	NIST	METEOR	BLEU	NIST	METEOR
Baseline	0.2009	5.866	0.2719	0.2051	5.957	0.2782	0.2003	5.889	0.2720
Context	<b>0.2039</b>	<b>5.941</b>	<b>0.2784</b>	<b>0.2088</b>	<b>6.036</b>	<b>0.2826</b>	0.2016	<b>5.956</b>	<b>0.2772</b>

Table 5: English  $\rightarrow$  German shared task system results using WMT-08 Europarl parallel data for training, dev06 for tuning, and three test sets, including the final 2008 test set. The row labeled “Context” uses the top-performing feature set {2 POS on left, 1 word on right}. Boldface marks scores that are significantly higher than the baseline.

models, including the lexicalized reordering model and the lexical translation model in the Moses MT system, or hierarchical or syntactic models (Chiang, 2005). Additional linguistic analysis (e.g., morphological disambiguation, named entity recognition, semantic role labeling) can be used to define new context features.

## 8 Conclusion

We have described a straightforward, scalable method for improving phrase translation models by modeling features of a phrase’s source-side context. Our method allows incorporation of features from any kind of source-side annotation and barely affects the decoding algorithm. Experiments show performance rivaling or exceeding strong, state-of-the-art baselines on standard translation tasks. Automatic feature selection can be used to achieve performance gains with just two or three context features. Performance is strongest when large in-domain training sets and high-accuracy NLP tools for the source language are available.

## Acknowledgments

This research was supported in part by an ARCS award to the first author, NSF IIS-0713265, supercomputing resources provided by Yahoo, and a Google grant. We thank Abhaya Agarwal, Ashish Venugopal, and Andreas Zollmann for helpful conversations and Joy Zhang for his Chinese segmenter. We also thank the anonymous reviewers for helpful comments.

## A Dataset Details (Chinese-English)

We trained on data from the NIST MT 2008 constrained Chinese-English track: Hong Kong Hansards and news (LDC2004T08), Sino-rama (LDC2005T10), FBIS (LDC2003E14), Xinhua (LDC2002E18), and financial news

(LDC2006E26)—total 2.5M sents., 66M Chinese words, 68M English. For news experiments, the newswire portion of the NIST 2004 test set was used for tuning, the full NIST 2003 test set was used for developmental testing and feature selection, and the NIST 2005 test set was used for testing (900-1000 sents. each). We also used the United Nations parallel text (LDC2004E12), divided into training (4.7M sents.; words: 136M Chinese, 144M English), tuning (2K sents.), and test sets (2K sents.). We removed sentence pairs where either side was longer than 80 words, segmented all Chinese text automatically,<sup>8</sup> and parsed/tagged using the Stanford parser with the pre-trained “xinhuaPCFG” model (Klein and Manning, 2003). Trigram language models were trained on the English side of the parallel corpus along with approximately 115M words from the Xinhua section of the English Gigaword corpus (LDC2005T12), years 1995–2000 (total 326M words).

## B Dataset Details (German-English)

For German  $\leftrightarrow$  English experiments, we used data provided for the WMT-07 shared task (1.1M sents., 26M German words, 27M English). We used dev06 for tuning, devtest06 for feature selection and developmental testing, and test07 for final testing (2K sents. each). We removed sentence pairs where either side was longer than 50 words and parsed/tagged the German and English data using the Stanford parser (Klein and Manning, 2003) (with pre-trained “germanFactored” and “englishPCFG” models). 5-gram language models were trained on the entire target side of the parallel corpus (37M German words, 38M English).

<sup>8</sup>Available at <http://projectile.is.cs.cmu.edu/research/public/tools/segmentation/lrsegmenter/lrSegmenter.perl>.



## References

- S. Banerjee and A. Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proc. of ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for MT and/or Summarization*.
- D. M. Bikel. 2004. A distributional analysis of a lexicalized statistical parsing model. In *Proc. of EMNLP*.
- T. Brants, A. C. Popat, P. Xu, F. J. Och, and J. Dean. 2007. Large language models in machine translation. In *Proc. of EMNLP-CoNLL*.
- P. F. Brown, J. Cocke, S. A. Della Pietra, V. J. Della Pietra, F. Jelinek, J. D. Lafferty, R. L. Mercer, and P. S. Roossin. 1990. A statistical approach to machine translation. *Computational Linguistics*, 16(2):79–85.
- C. Callison-Burch, P. Koehn, C. Fordyce, and C. Monz, editors. 2007. *Proc. of the 2nd Workshop on Statistical Machine Translation*.
- M. Carl and A. Way. 2003. *Recent Advances in Example-Based Machine Translation*. Kluwer Academic.
- M. Carpuat and D. Wu. 2007. Improving statistical machine translation using word sense disambiguation. In *Proc. of EMNLP-CoNLL*.
- Y. Chan, H. Ng, and D. Chiang. 2007. Word sense disambiguation improves statistical machine translation. In *Proc. of ACL*.
- S. Chen and J. Goodman. 1998. An empirical study of smoothing techniques for language modeling. Technical report 10-98, Harvard University.
- D. Chiang. 2005. A hierarchical phrase-based model for statistical machine translation. In *Proc. of ACL*.
- G. Doddington. 2002. Automatic evaluation of machine translation quality using  $n$ -gram co-occurrence statistics. In *Proc. of HLT*.
- D. Gildea. 2001. Corpus variation and parser performance. In *Proc. of EMNLP*.
- J. Giménez and L. Màrquez. 2007. Context-aware discriminative phrase selection for statistical machine translation. In *Proc. of the 2nd Workshop on Statistical Machine Translation*.
- I. Guyon and A. Elisseeff. 2003. An introduction to variable and feature selection. *Journal of Machine Learning Research*, 3:1157–1182.
- D. Klein and C. D. Manning. 2003. Fast exact inference with a factored model for natural language parsing. In *Advances in NIPS 15*.
- P. Koehn, F. J. Och, and D. Marcu. 2003. Statistical phrase-based translation. In *Proc. of HLT-NAACL*, pages 127–133.
- P. Koehn, H. Hoang, A. Birch, C. Callison-Burch, M. Federico, N. Bertoldi, B. Cowan, W. Shen, C. Moran, R. Zens, C. Dyer, O. Bojar, A. Constantin, and E. Herbst. 2007. Moses: Open source toolkit for statistical machine translation. In *ACL demonstration session*.
- P. Koehn. 2004. Statistical significance tests for machine translation evaluation. In *Proc. of EMNLP*.
- José B. Marino, Rafael E. Banchs, Josep M. Crego, Adrià de Gispert, Patrik Lambert, José A. R. Fonollosa, and Marta R. Costa-jussà. 2006.  $N$ -gram-based machine translation. *Computational Linguistics*, 32(4):527–549.
- M. Nagao. 1984. A framework of a mechanical translation between Japanese and English by analogy principle. In *Proc. of the International NATO Symposium on Artificial and Human Intelligence*. Elsevier North-Holland, Inc.
- NIST. 2006. NIST 2006 machine translation evaluation official results.
- F. Och and H. Ney. 2003. A systematic comparison of various statistical alignment models. *Computational Linguistics*, 29(1).
- F. J. Och. 2003. Minimum error rate training for statistical machine translation. In *Proc. of ACL*, pages 160–167.
- K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu. 2002. BLEU: a method for automatic evaluation of machine translation. In *Proc. of ACL*.
- H. Somers. 1999. Review article: Example-based machine translation. *Machine Translation*, 14(2).
- A. Stolcke. 2002. SRILM—an extensible language modeling toolkit. In *Proc. of ICSLP*.
- N. Stroppa, A. van den Bosch, and A. Way. 2007. Exploiting source similarity for SMT using context-informed features. In *Proc. of TMI*.
- D. Vickrey, L. Biewald, M. Teyssier, and D. Koller. 2005. Word-sense disambiguation for machine translation. In *Proc. of HLT-EMNLP*.