

# Constraining the Phrase-Based, Joint Probability Statistical Translation Model

Alexandra Birch Chris Callison-Burch Miles Osborne

School of Informatics  
University of Edinburgh  
Buccleuch Place  
Edinburgh, EH8 9LW  
a.c.birch-mayne@sms.ed.ac.uk

## Abstract

The joint probability model proposed by Marcu and Wong (2002) provides a strong probabilistic framework for phrase-based statistical machine translation (SMT). The model's usefulness is, however, limited by the computational complexity of estimating parameters at the phrase level. We present a method of constraining the search space of the joint probability model based on statistically and linguistically motivated word alignments. This method reduces the complexity and size of the joint model and allows it to display performance superior to the standard phrase-based models.

## 1 Introduction

Machine translation is a hard problem because of the highly complex, irregular and diverse nature of natural languages. It is impossible to accurately model all the linguistic rules that shape the translation process, and therefore a principled approach uses statistical methods to make optimal decisions given incomplete data.

The original IBM Models (Brown et al., 1993) learned only word-to-word alignment probabilities which made it computationally feasible to estimate model parameters from large amounts of training data. Phrase-based SMT models, such as the alignment template model (Och, 2003), improve on word-based models because phrases provide local context which leads to better lexical choice and more reliable local reordering. However, most phrase-based models extract their phrase pairs from previously word-aligned corpora using ad-hoc heuristics. These models perform no search for optimal phrasal alignments.

Even though this is an efficient strategy, it is a departure from the rigorous statistical framework of the IBM Models.

Marcu and Wong (2002) proposed the joint probability model which directly estimates the phrase translation probabilities from the corpus in a theoretically governed way. This model neither relies on potentially sub-optimal word alignments nor on heuristics for phrase extraction. Instead, it searches the phrasal alignment space, simultaneously learning translation lexicons for both words and phrases. The joint model has been shown to outperform standard models on restricted data sets such as the small data track for Chinese-English in the 2004 NIST MT Evaluation (Przybocki, 2004).

However, considering all possible phrases and all their possible alignments vastly increases the computational complexity of the joint model when compared to its word-based counterpart. This results in prohibitively slow training and heavy use of memory resources. The large size of the model means that only a very small proportion of the alignment space can be searched, and this reduces the chances of finding optimum parameters. Furthermore, the complexity of the joint model makes it impossible to scale up to the larger training corpora available today, preventing the model from being more widely adopted.

In this paper, we propose a method of constraining the search space of the joint model to areas where most of the unpromising phrasal alignments are eliminated and yet as many potentially useful alignments as possible are still explored. The joint model is constrained to phrasal alignments which do not contradict a set high confidence word alignments for each sentence. These high confidence alignments can incorporate information from both statistical and linguistic sources. We show that by

using the points of high confidence from the intersection of the bi-directional Viterbi alignments to reduce complexity, translation quality also improves. We also show that the addition of linguistic information from a machine readable dictionary and aligning identical words further improves the model.

Apart from the large memory requirements of the joint model, it is computationally very expensive to train. We describe a modification to the Expectation Maximisation (EM) algorithm which greatly increases the speed of the training without compromising the quality of the resulting translations.

## 2 Models

### 2.1 Standard Phrase-based Model

Most phrase-based models (Och, 2003; Koehn et al., 2003; Vogel et al., 2003) rely on a pre-existing set of word-based alignments from which they induce their parameters. In this project we use the model described by Koehn et al. (2003) which extracts its phrase alignments from a corpus that has been word aligned. From now on we refer to this phrase-based model as the standard model.

The standard model decomposes the foreign input sentence  $F$  into a sequence of  $I$  phrases  $\bar{f}_1, \dots, \bar{f}_I$ . All segmentations are assumed to be equally probable. Each foreign phrase  $\bar{f}_i$  is translated to an English phrase  $\bar{e}_i$  using the probability distribution  $\theta(\bar{f}_i|\bar{e}_i)$ . English phrases may be re-ordered using a relative distortion probability  $d(\cdot)$ . The model is defined as follows:

$$p(F|E) = \prod_{i=1}^I \theta(\bar{f}_i|\bar{e}_i)d(\cdot) \quad (1)$$

As alignments between phrases are constructed from word alignments, there is no summing over possible alignments. This model performs no search for optimal phrase pairs. Instead, it extracts phrase pairs  $(\bar{f}_i, \bar{e}_i)$  in the following manner. First, it uses the IBM Models to learn the Viterbi alignments for English to Foreign and Foreign to English. It then uses a heuristic to reconcile the two alignments, starting from the points of high confidence in the intersection of the two Viterbi alignments and growing towards the points in the union. Points from the union are selected if they are adjacent to points from the intersection and their words are previously unaligned. Koehn

et al. (2003) discusses and compares variations on this strategy.

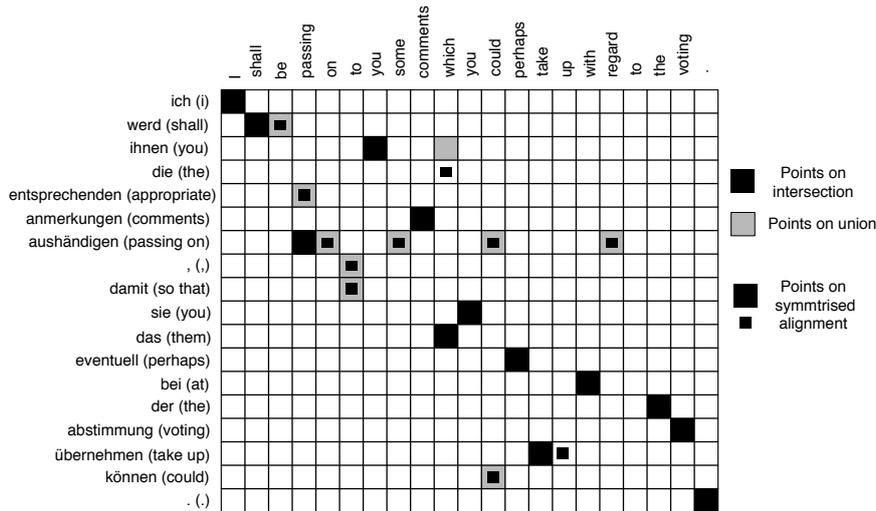
Phrases are then extracted by selecting phrase pairs which are ‘consistent’ with the symmetrised alignment. Here ‘consistent’ means that all words within the source language phrase are only aligned to the words of the target language phrase and vice versa. Finally the phrase translation probability distribution is estimated using the relative frequencies of the extracted phrase pairs.

This approach to phrase extraction means that phrasal alignments are locked into the symmetrised alignment. This is problematic for two reasons: firstly, the symmetrisation process will grow an alignment based on arbitrary decisions about adjacent words, and secondly, because word alignments inadequately represent the real dependencies between translations.

In Figure 1 we can see an example of symmetrization from the German-English Europarl corpus when training the standard model on nearly 500,000 sentences.

Figure 1 highlights the problems of modelling the translation process using word alignments. Parallel sentences are often not literally translated: some words are dropped, like ‘entsprechenden (appropriate)’, and some have been translated in a way for which there is no word alignment possible, like the words in ‘with regard to’. The points from the intersection are mostly correct, and the symmetrization process adds some new points correctly. For instance, the word ‘on’ correctly gets added to alignment with ‘aushändigen’ to create the phrase ‘passing on’. However, further words get incorrectly aligned because they are adjacent to points on the symmetrized alignment, such as ‘, damit’ with ‘to’.

Locking the set of phrase pairs into the symmetrized alignment is clearly not ideal. By heuristically creating phrasal alignments from Viterbi word-level alignments, we throw away a great deal of the information that was estimated when learning word alignment parameters and we can introduce errors. With the joint model one can search areas of the alignment space. Common phrases such as ‘with regard to’ are treated as a unit, improving its chances of finding good parameters. This allows us to learn a distribution of possible phrasal alignments that better handles the uncertainty inherent in the translation process.



**Figure 1.** Example of symmetrisation which illustrates that errors can be introduced through heuristic merging, e.g. ‘, damit’ is aligned with ‘to’

## 2.2 Joint Probability Model

The joint probability model (Marcu and Wong, 2002), does not rely on a pre-existing set of word-level alignments. Like the IBM Models, it uses Expectation Maximisation to align and estimate the probabilities for sub-sentential units in a parallel corpus. Unlike the IBM Models, it does not constrain the alignments to being single words.

The basic model is defined as follows. Phrases are created from words and commonly occurring sequences of words. Concepts,  $c_j$ , are defined as a pair of aligned phrases  $\langle \bar{e}_i, \bar{f}_i \rangle$ . A set of concepts which completely covers the sentence pair is denoted by  $C$ . Phrases are restricted to being sequences of words which occur above a certain frequency in the corpus. Commonly occurring phrases are more likely to lead to the creation of useful phrase pairs, and without this restriction the search space would be much larger.

The probability of a sentence and its translation is the sum of all possible alignments,  $C$  each of which is defined as the product of the probability of all individual concepts:

$$p(F, E) = \sum_{C \in \mathcal{C}} \prod_{\langle \bar{e}_i, \bar{f}_i \rangle \in C} p(\langle \bar{e}_i, \bar{f}_i \rangle) \quad (2)$$

The model is trained by initialising the translation table and then performing EM as described below.

### 2.2.1 Initialising Translation Table

Before starting EM all phrasal alignments are assumed to be equally probable. Under these circumstances, the probability of a concept  $c_j$  in sentences  $(E, F)$  is equal to the number of phrasal alignments which contain this concept divided by the total number of phrasal alignments that can be built between the two sentences. This probability can be approximated by using the lengths of the two phrases and the lengths of the two sentences with Stirling numbers of the second kind as described by Marcu and Wong (2002). We are thus able to initialise all possible alignments.

The size of the translation table is largely determined by the initialisation phase, and so it greatly impacts on the scalability of the model.

### 2.2.2 Expectation Maximisation

After initialising the translation parameters, alignments will have different probabilities. It is no longer possible to collect fractional counts over all possible alignments in polynomial time. EM is therefore performed approximately to improve parameters and increase the probability of the corpus.

An iteration of EM starts by creating an initial phrasal alignment of high probability. This is done by selecting the highest probability concepts that cover the sentence pair. Then the model hill-climbs towards the optimal Viterbi alignment by using a set of modifying operations. These operations break and merge concepts, swap words between concepts and move words across concepts.

The model calculates the probabilities associated with all alignments generated in this process and collects fractional counts for the concepts based on these probabilities.

### 2.2.3 Complexity

Training the IBM models is computationally challenging, but the joint model is much more demanding. Considering all possible segmentations of phrases and all their possible alignments vastly increases the number of possible alignments that can be formed between two sentences.

<i>E</i> Length	<i>F</i> Length	No. Alignments
5	5	6721
10	10	818288740923
20	20	4.4145633531e+32
40	40	2.7340255177e+83

**Table 1.** The number of possible phrasal alignments for sentence pairs calculated using Stirling numbers of the second kind.

Table 1 shows just how many phrasal alignments are possible between sentences of different length. Even for medium length sentences that are 20 words in lengths, the total number of alignments is huge. Apart from being intractable, when one has a very large parameter estimation space the EM algorithm struggles to discover good parameters. Pereira and Schabes (1992) proposed a method for dealing with this problem for PCFG estimation from treebanks. They encouraged the probabilities into good regions of the parameter space by constraining the search to only consider parses that did not cross Penn-Treebank nodes.

## 3 Constraining the Joint Model

The Joint Model requires a strategy for restricting the search for phrasal alignments to areas of the alignment space which contain most of the probability mass. We propose a method which examines phrase pairs that are consistent with the set of high confidence word alignments defined for the sentence. By ‘consistent’ we mean that for a concept  $\langle \bar{e}_i, \bar{f}_i \rangle$  to be valid, we make sure that if any word in  $\bar{e}_i$  is part of a high confidence alignment, then the word to which it is aligned must be included in  $\bar{f}_i$  and vice versa. Phrases must still occur above a certain frequency in the corpus to be considered.

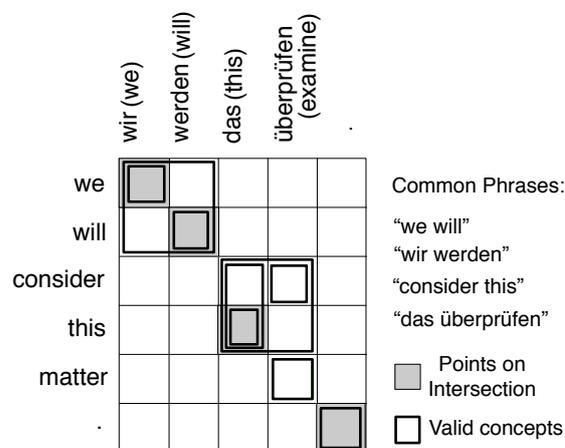
The constraints on the model are applied during the initialisation phase of the training. During EM,

it would be very computationally expensive to enforce constraints and there would be less benefit as only high probability alignments are visited.

### 3.1 IBM Constraints

The standard phrase-based model is based on a complex series of models, parameters and heuristics which allow it to be efficient. The joint probability model is a more principled and conceptually simpler model but it is very inefficient. By using the IBM Models to constrain the joint model, we are searching areas in the phrasal alignment space where both models overlap. We combine the advantage of prior knowledge about likely word alignments with the ability to perform a probabilistic search around them.

The IBM constraints are the high confidence word alignments that result from taking the intersection of the bi-directional Viterbi alignments. This strategy for extracting phrase pairs that are coherent with the Viterbi alignments is similar to that of the standard phrase-based model. However, the constrained joint model does not lock the search into a heuristically derived symmetrized alignment.



**Figure 2.** The area of alignment space searched using IBM constraints for an example sentence.

Figure 2 shows us the space searched by the model for an example sentence. All valid concepts are consistent with all high confidence word alignments and either comprise of words or commonly occurring phrases. The concept  $\langle \text{‘wir’, ‘consider’} \rangle$  would break the high confidence alignment between ‘wir’ and ‘we’ and would therefore be invalid. We can also see that the model searches more intensively areas

of the sentence about which there is little certainty. Searching over an area of lower probability is preferable to using a heuristic to arbitrarily align all unaligned words. Searching allows good phrasal alignments to be discovered, for instance <‘das überprüfen’, ‘consider this’ >.

### 3.2 Linguistic Constraints

By constraining the joint model using high confidence word alignments, any external knowledge sources can be included into the probabilistic framework. Linguistic constraints can be combined to guide the training of the joint model. In this paper we use a bilingual dictionary and identical words to contribute further alignment points. These constraints are combined in a simple linear fashion. First the IBM constraints are collected, then identical words that do not contradict the IBM alignments are aligned. Finally, word entries from the dictionary are used to align as yet unaligned words.

These linguistic constraints are useful with small sets of training data, but for larger corpora, dictionaries and identical words would contribute less to the quality of the final translations. However, the advantage of being able to include any knowledge about word alignments within a statistical model is compelling.

## 4 Optimizing the Joint Model

### 4.1 Prior counts from Word-Aligned Corpora

The joint probability model can only be trained with small amounts of parallel data and consequently the resulting parameters suffer from sparse counts. In order to make fractional counts more reliable, we can include information which encodes our prior belief about word-to-word alignments. This is desirable as word alignments are less prone to sparse statistics than phrasal alignments.

When training the joint model, we have initially assumed a uniform probability across all possible alignments. In a sentence, concepts of the same size will be assigned the same fractional counts. If one concept occurs more often over the entire corpus, its final parameter value will be higher. However, when the training corpus is very small, it is unlikely for the model to have seen representative occurrences of the concepts.

In order to overcome this problem, the joint model can use information about word-alignments

generated by the IBM models. A simple way to include this knowledge is to use the high confidence points from the intersection of the bi-directional Viterbi alignments. Concepts which contain many points of high confidence will be more probable than concepts of the same size which contain none.

We define a prior count which reflects the probability of the phrasal alignment given the high confidence word alignments:

$$pc(\bar{e}, \bar{f}) = \frac{|align|}{\min(|\bar{e}|, |\bar{f}|)}$$

We divide the number of word alignments contained within the concept by the total number of possible word alignments for the concept, which is equal to the length of the shorter of the two phrases. We add a small fraction (0.1) to both the numerator and the denominator to smooth and avoid zero probabilities.

One way to include this prior count in the model would be to calculate it separately and then use it in the decoding process as one of the features of the log linear model. This would be similar to the lexical weighting employed by Koehn et al. (2003). In the joint model, however, we must perform EM and including these probabilities in the training of the model will improve the overall quality of alignments searched. These counts are thus included in the initialisation phase of the joint model training with the calculation of the fractional counts:

$$fc(\bar{e}, \bar{f}) = (1 - \lambda)p(\bar{e}, \bar{f}|E, F) + \lambda pc(\bar{e}, \bar{f})$$

The fractional count for each concept in each sentence is calculated by interpolating the joint probability of the concept, based on the Stirling numbers, and the prior count, which reflects the probability of the phrasal alignment given the high confidence word alignments. The use of the weight to balance the two contributions allows us to adjust for differences in scale and our confidence in each of the two measures. After testing various settings for  $\lambda$  the value 0.5 gave the best Bleu scores. Callison-Burch et al. (2004) used a similar technique for combining word and sentence aligned data. However, they inserted data from labelled word alignments which meant that they did not need to sum over all possible alignments for a sentence pair.

## 4.2 Fast Hill-climbing

The constraints on the joint model reduce its size by restricting the initialisation phase of the training. This is one of the two major drawbacks of the model discussed by Marcu and Wong (2002). The other major drawback is the computational cost of the training procedure. Fast hill-climbing is necessary to make EM training more tractable.

The joint model examines all possible swaps, splits, merges and moves for the set of concepts that have been selected as part of the initial alignment. Normal hill-climbing repeatedly performs a very expensive search over all possible steps, selecting the best step each time and applying it until no further improvement is found. In fast hill-climbing, instead of selecting only the best step, we collect all the steps that improve the probability of the initial phrasal alignment, and only search once. We then apply them one by one to the initial phrasal alignment.

This approach has the disadvantage of heavily weighting the initial alignment. All alignments generated during one iteration of EM are only one step away from the initial alignment, so the counts for the concepts in this alignment will be high. This is a drastic change to Viterbi training, but such measures are needed to reduce the training time from nearly 5 hours to complete one iteration of EM for just 5000 sentences.

## 5 Experiments

We used the German-English Europarl corpus (Koehn, 2002) to perform our experiments. Europarl contains proceedings from the European Parliament covering the years 1996-2003. The test set consisted of 1755 sentences which ranged from 5 to 15 words in length.

For the language model we used the SRI Language Modelling Toolkit (Stolcke, 2002) to train a trigram model on the English section of the Europarl corpus.

The first baseline used was the standard phrase-based model (Koehn et al., 2003) with the default feature set and options. The second baseline was an implementation of the joint probability model as described by Marcu and Wong (2002). We set all the model's maximum phrase length to four words long.

To perform the translations we used the Pharaoh (Koehn, 2004) beam search decoder version 1.2.8, with all the standard settings. Our

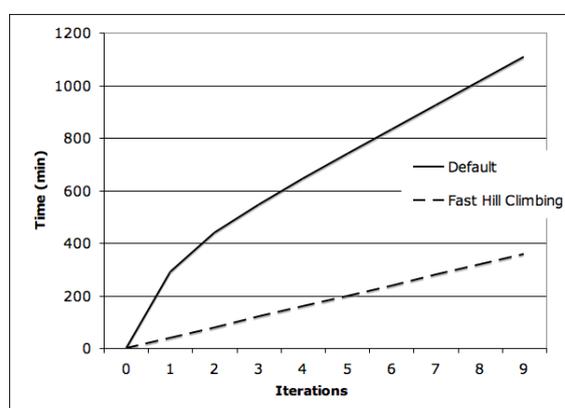
evaluation metric was Bleu (Papineni et al., 2002) which compares the output sentences with human translated sentences using 4-gram precision.

We also perform experiments with a bilingual dictionary which comes with Ding, an open source translation program<sup>1</sup>.

## 6 Results

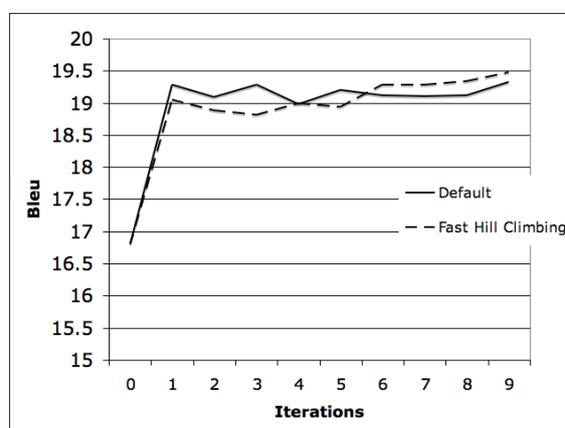
### 6.1 Fast Hill-climbing

EM training of the joint model was prohibitively slow even for the smallest data sets, so the first experiment explores the gains to be made by using fast hill-climbing.



**Figure 3.** Time taken for EM training in minutes per iteration for 5,000 sentences on a machine with 2Gb RAM and a 2.4GHz CPU

In Figure 3 we can see that fast hill-climbing is much faster than the normal hill-climbing. We have reduced the time taken to perform the first iteration from nearly 5 hours to about 40 minutes, which is about a factor of eight.



**Figure 4.** Bleu scores using 5,000 sentences training data

<sup>1</sup>See <http://www-user.tu-chemnitz.de/~fri/ding/>

The effect of fast hill-climbing on the quality of translations can be seen in Figure 4. The default method slightly outperforms fast hill-climbing for the first few iterations, but then fast hill-climbing overtakes it. The difference in performance between the two methods is small and we apply fast hill-climbing in the remaining experiments.

## 6.2 IBM Constraints

Corpus Size	10,000	20,000	40,000
Standard Model	21.69	23.61	25.52
Joint Model	19.93	-	-
+ IBM	22.13	23.08	24.16
+ IBM + prior	22.79	24.33	25.99

**Table 2.** Bleu scores for the joint model with IBM constraints and prior counts, corpus size indicates number of sentence pairs

Corpus Size	10,000	20,000	40,000
Standard Model	95k	200k	405k
Joint Model	6,178k	-	-
+ IBM	1,457k	2,738k	4,993k
+ IBM + prior	1,451k	2,724k	4,964k

**Table 3.** Translation table size in number of phrase pairs

In Tables 2 and 3 we can see the differences in size and performance between the baseline model and the joint model for different sizes of training corpora. The unconstrained joint model produces a very large translation table, containing more than 6 million phrase pairs. The size of the model hampers its performance, resulting in a poor Bleu score. In fact it was only able to be trained on a maximum of 10,000 sentences before running out of memory on a machine with 2Gb of RAM.

By using IBM constraints, the performance of the joint model improves, beating even the standard phrase-based model. The resulting translation table is, however, still quite large and only about four times smaller than the unconstrained joint model. On examining the phrase pairs produced for each sentence, we discovered that the reason for the large size of the model was due to longer sentences for which there were few points of high confidence.

Table 2 shows that adding prior counts based on word alignments to the initial estimation of the joint probability improves the Bleu score.

## 6.3 Linguistic Constraints

The effect of adding linguistic constraints to the IBM word constraints is shown in tables 4 and 5.

Corpus Size	10,000	20,000	40,000
Standard Model	95k	200k	405k
Joint + IBM + prior	1,451k	2,724k	4,964k
+ Ident.	1,361k	2,557k	4,649k
+ Ident. + Dict.	1,096k	2,079k	3,834k

**Table 4.** Translation table size in number of phrase pairs when linguistic constraints are added to the joint model

Corpus Size	10,000	20,000	40,000
Standard Model	21.69	23.61	25.52
Joint + IBM + prior	22.79	24.33	25.99
+ Ident.	23.30	24.90	26.12
+ Ident. + Dict.	23.20	24.96	26.13

**Table 5.** Bleu scores for different training corpus sizes

Table 4 shows that by adding high confidence alignments for identical words and forcing phrase pairs to be consistent with these as well as the IBM constraints, we reduce the size of the model but only slightly. Including points from the bilingual dictionary results in a sizeable reduction of about 20%.

Table 5 shows that the inclusion of lexical information into the model improves performance. The improvement in Bleu score seems to reduce with the increase in training data. As the model is trained on more data, external knowledge sources provide less advantage.

## 7 Conclusion

In this paper we have shown that using the joint probability model to estimate phrase translation probabilities results in a better performance than the standard heuristic approach. This suggests that there are gains to be had by using a more principled statistical framework.

We presented the first attempt at constraining the joint probability model. By introducing constraints to the alignment space we can greatly reduce the complexity of the model and increase its performance. The strategy of using IBM constraints with the joint model allows it to search areas of the alignment space with a higher probability mass, resulting in better parameters. A constrained joint probability model can train on larger

corpora making the model more widely applicable. Also, our particular method of constraining the joint model makes it easy to include linguistic information into the probabilistic framework of SMT.

When using IBM constraints, the search space of the joint model is comparable to that of standard phrase-based models, but does not depend on ad-hoc heuristics for phrase pair extraction. If the joint model were to be further engineered for efficiency, for example by using word classes, it could potentially replace the standard models for smaller data conditions.

## References

- Peter F. Brown, Stephen A. Della Pietra, Vincent J. Della Pietra, and Robert L. Mercer. 1993. The mathematics of machine translation: Parameter estimation. *Computational Linguistics*, 19(2):263–311.
- Chris Callison-Burch, David Talbot, and Miles Osborne. 2004. Statistical machine translation with word- and sentence-aligned parallel corpora. In *Proceedings of ACL*.
- Philipp Koehn, Franz Och, and Daniel Marcu. 2003. Statistical phrase-based translation. In *Proceedings of HLT/NAACL*, pages 127–133.
- Philipp Koehn. 2002. Europarl: A multilingual corpus for evaluation of machine translation. Unpublished.
- Philipp Koehn. 2004. Pharaoh: a beam search decoder for phrase-based statistical machine translation models. In *Proceedings of AMTA*, pages 115–124.
- Daniel Marcu and William Wong. 2002. A phrase-based, joint probability model for statistical machine translation. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing, EMNLP*.
- Franz Josef Och and Hermann Ney. 2000. Improved statistical alignment models. In *Proceedings of ACL*, pages 440–447.
- Franz Josef Och. 2003. *Statistical Machine Translation: From Single-Word Models to Alignment Templates*. Ph.D. thesis, RWTH Aachen Department of Computer Science, Aachen, Germany.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: a method for automatic evaluation of machine translation. In *Proceedings of ACL*, pages 311–318.
- Fernando Pereira and Yves Schabes. 1992. Inside-outside reestimation from partially bracketed corpora. In *Proceedings of ACL*, pages 128–135.
- Mark Przybocki. 2004. NIST 2004 machine translation evaluation results. Confidential e-mail to workshop participants, May.
- Andreas Stolcke. 2002. SRILM - an extensible language modeling toolkit. In *Proceedings of Spoken Language Processing*, pages 901–904.
- Stephan Vogel, Ying Zhang, Fei Huang, Alicia Tribble, Ashish Venugopal, Bing Zhao, and Alex Waibel. 2003. The CMU statistical machine translation system. In *Machine Translation Summit*.