

Exact Decoding with Multi Bottom-Up Tree Transducers*

Daniel Quernheim

Institut für Maschinelle Sprachverarbeitung

University of Stuttgart

daniel@ims.uni-stuttgart.de

Abstract

We present an experimental statistical tree-to-tree machine translation system based on the multi-bottom up tree transducer including rule extraction, tuning and decoding. Thanks to input parse forests and a “no pruning” strategy during decoding, the obtained translations are competitive. The drawbacks are a restricted coverage of 70% on test data, in part due to exact input parse tree matching, and a relatively high runtime. Advantages include easy redecoding with a different weight vector, since the full translation forests can be stored after the first decoding pass.

1 Introduction

In this contribution, we present an implementation of a translation model that is based on ℓ -XMBOT (the multi bottom-up tree transducer of Arnold and Dauchet (1982) and Lilin (1978)).¹ Intuitively, an MBOT is a synchronous tree sequence substitution grammar (STSSG, Zhang et al. (2008a); Zhang et al. (2008b); Sun et al. (2009)) that has discontinuities only on the target side (Maletti, 2011). From an algorithmic point of view, this makes the MBOT more appealing than STSSG as demonstrated by Maletti (2010). Formally, MBOT is expressive enough to express all sensible translations (Maletti, 2012)². Figure 2 displays sample rules of the MBOT variant, called ℓ -XMBOT,

This work was supported by Deutsche Forschungsgemeinschaft grant MA/4959/1–1.

¹The system presented in this paper is variant of the system presented at last year’s workshop (Quernheim and Cap, 2014), without morphological enhancements.

²A translation is sensible if it is of linear size increase and can be computed by some (potentially copying) top-down tree transducer.

that we use (in a graphical representation of the trees and the alignment). Recently, a shallow version of MBOT has been integrated into the popular Moses toolkit (Braune et al., 2013). Our implementation is exact in the sense that it does absolutely no pruning during decoding and thus preserves all translation candidates, while having no mechanism to handle unknown structures. (We added dummy rules that leave unseen lexical material untranslated.) The coverage is thus limited, but still considerably high. Source-side and target-side syntax restrict the search space so that decoding stays tractable. Only the language model scoring is implemented as a separate reranker. This has several advantages: (1) We can use input parse forests (Liu et al., 2009). (2) Not only is the output optimal with regard to the theoretical model, also the space of translation candidates can be efficiently stored as a weighted regular tree grammar. The best translations can then be extracted using the k-best algorithm by Huang and Chiang (2005). Rule weights can be changed without the need for explicit redecoding, the parameters of the log-linear model can be changed, and even new features can be added. These properties are especially helpful in tuning, where only the k-best algorithm has to be re-run in each iteration. A model in similar spirit has been described by Huang et al. (2006); however, it used target syntax only (using a top-down tree-to-string transducer backwards), and was restricted to sentences of length at most 25. We do not make such restrictions.

The theoretical aspects of ℓ -XMBOT and their use in our translation model are presented in Section 2. Based on this, we implemented a machine translation system that we are going to make available to the public. Section 4 presents the most important components of our ℓ -XMBOT implemen-

tation, and Section 5 presents our submission to the WMT15 shared translation task.

2 Theoretical Model

In this section, we present the theoretical generative model that is used in our approach to syntax-based machine translation: the multi bottom-up tree transducer (Maletti, 2011). It is a variant of the linear and nondeleting extended multi bottom-up tree transducers without states. We omit the technical details and give graphical examples only to illustrate how the device works, but refer to the literature for the theoretical background. Roughly speaking, a local multi bottom-up tree transducer (ℓ MBOT) has rules that replace one nonterminal symbol N on the source side by a tree, and a sequence of nonterminal symbols on the target side linked to N by one tree each. These trees again have linked nonterminals, thus allowing further rule applications.

Our ℓ MBOT rules are obtained automatically from data like that in Figure 1. Thus, we (word) align the bilingual text and parse it in both the source and the target language. In this manner we obtain sentence pairs like the one shown in Figure 1. To these sentence pairs we apply the rule extraction method of Maletti (2011). The rules extracted from the sentence pair of Figure 1 are shown in Figure 2. Note the discontinuous alignment of *went* to *ist* and *gegangen*, resulting in discontinuous rules.

The application of those rules is illustrated in Figure 3 (a *pre-translation* is a pair consisting of a source tree and a sequence of target trees). While it shows a synchronous derivation, our main use case of ℓ MBOT rules is *forward application* or *input restriction*, that is the calculation of all target trees that can be derived given a source tree. For a given synchronous derivation d , the source tree generated by d is $s(d)$, and the target tree is $t(d)$. The yield of a tree is the string obtained by concatenating its leaves.

The theoretical justification for decomposing the translation model into a source model and a target model is a theorem that states that every ℓ MBOT can be replaced by a composition of a linear nondeleting extended top-down tree transducer (XTOP) and a linear homomorphic MBOT (Engelfriet et al., 2009). We implemented the first step of the composition as an XTOP that generates possible derivation trees. States in this de-

vice are linked nonterminals in the ℓ MBOT rules, and it translates left-hand sides into rule identifiers. The second step is implemented as a homomorphic multi bottom-up tree transducer. While we construct the first step of the composition explicitly, we only use the second device to evaluate single trees.

Apart from ℓ MBOT application to input trees, we can even apply ℓ MBOT to *parse forests* and even *weighted regular tree grammars* (RTGs) (Fülöp and Vogler, 2009). RTGs offer an efficient representation of weighted forests, which are sets of trees such that each individual tree is equipped with a weight. This representation is even more efficient than packed forests (Mi et al., 2008) and moreover can represent an infinite number of weighted trees. The most important property that we utilize is that the output tree language is regular, so we can represent it by an RTG (cf. preservation of regularity (Maletti, 2011)). Indeed, every input tree can only be transformed into finitely many output trees by our model, so for a given finite input forest (which the output of the parser is) the computed output forest will also be finite and thus regular.

3 Translation Model

Given a source language sentence e and corresponding weighted parse forest $F(e)$, our translation model aims to find the best corresponding target language translation \hat{g} ;³ i.e.,

$$\hat{g} = \arg \max_g p(g|e) .$$

We estimate the probability $p(g|e)$ through a logarithmic combination of component models with parameters λ_m scored on the derivations d such that the source tree $s(d)$ of d is in the parse forest of e and the yield of the target tree $t(d)$ reads g . With

$$D(e, g) = \{d \mid s(d) \in F(e) \text{ and } \text{yield}(t(d)) = g\},$$

we thus have:⁴

$$p(g|e) \propto \sum_{d \in D(e, g)} \prod_{m=1}^{11} h_m(d)^{\lambda_m}$$

Our model uses the following features $h_m(\cdot)$ for a derivation:

³Our main translation direction is English to German.

⁴While this is the clean theoretical formulation, we make two approximations to $D(e, g)$: (1) The parser we use returns a pruned parse forest. (2) We only sum over derivations with the same target sentence that actually appear in the k-best list.

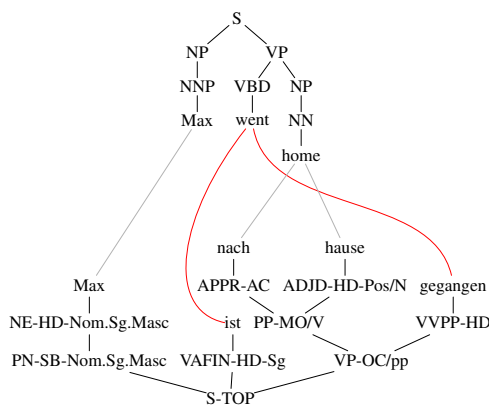


Figure 1: Aligned parsed sentences.

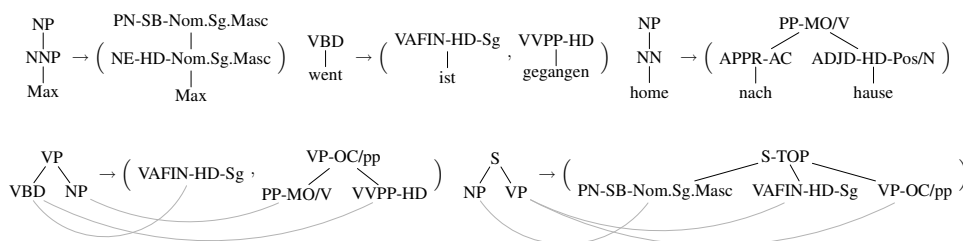


Figure 2: Extracted rules.

- (1) Translation weight normalized by source root symbol
- (2) Translation weight normalized by all root symbols
- (3) Lexical translation weight source \rightarrow target
- (4) Lexical translation weight target \rightarrow source
- (5) Target side language model: $p(g)$
- (6) Input parse tree probability assigned to $s(t)$ by the parser of e

The rule weights required for (1) are relative frequencies normalized over all extracted rules with the same root symbol on the left-hand side. In the same fashion the rule weights required for (2) are relative frequencies normalized over all rules with the same root symbols on both sides. The lexical weights for (3) and (4) are obtained by multiplying the word translations $w(g_i|e_j)$ [respectively, $w(e_j|g_i)$] of lexically aligned words (g_i, e_j) across (possibly discontinuous) target side sequences.⁵ Whenever a source word e_j is aligned to multiple target words, we average over the word

⁵The lexical alignments are different from the links used to link nonterminals.

translations:⁶

$$h_3(d) = \prod_{\substack{\text{lexical item} \\ e \text{ occurs in } s(d)}} \text{average} \{w(g|e) \mid g \text{ aligned to } e\}$$

4 Implementation

Our implementation is very close to the theoretical model and consists of several independent components, most of which are implemented in Python. The system does not have any dependencies other than the need for parsers for the source and target language, a word alignment tool and optionally an implementation of some tuning algorithm.

Rule extraction From a parallel corpus of which both halves have been parsed and word aligned, multi bottom-up tree transducer rules are extracted according to the procedure laid out in (Maletti, 2011). In order to handle unknown words, we add dummy identity translation rules for lexical material that was not present in the training data.

⁶If the word e_j has no alignment to a target word, then it is assumed to be aligned to a special NULL word and this alignment is scored.

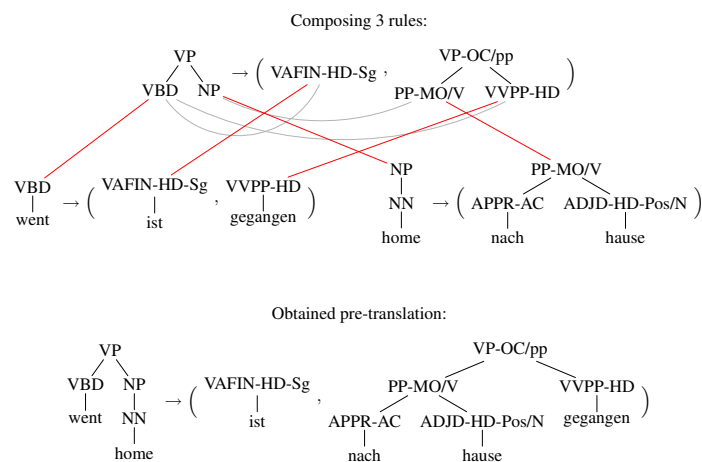


Figure 3: Synchronous rule application.

Translation model building Given a set of rules, translation weights (see above) are computed for each unique rule. The translation model is then converted into a source, a weight and a target model. The source model (an RTG represented in an efficient binary format) is used for decoding and maps input trees to trees over rule identifiers representing derivations. The weight model and the target model can be used to reconstruct the weight and the target realization of a given derivation.

Decoder For every input sentence, the decoder transforms a forest of parse trees to a forest of translation derivations by means of forward application. These derivations are trees over the set of rules (represented by rule identifiers). One of the most useful aspects of our model is the fact that decoding is completely independent of the weights, as no pruning is performed and all translation candidates are preserved in the translation forest. Thus, even after decoding, the weight model can be changed, augmented by new features, etc.; even the target model can be changed, e.g. to support parse tree output instead of string output. In all of our experiments, we used string output, but it is conceivable to use other realizations. For instance, a syntactic language model could be used for output tree scoring. Also, recasing is extremely easy when we have part-of-speech tags to base our decision on (proper names are typically uppercase, as are all nouns in German).

Another benefit of having a packed representation of all candidates is that we can easily check whether the reference translation is included in the candidate set (“force decoding”). The freedom to

allow arbitrary target models that rewrite derivations is related to current work on interpreted regular tree grammars (Koller and Kuhlmann, 2011), where arbitrary algebras can be used to compute a realization of the output tree.

k-best extractor From the translation derivation RTGs, a k-best list of derivations can be extracted (Huang and Chiang, 2005) very efficiently. This is the only step that has to be repeated if the rule weights or the parameters of the log-linear model change. The derivations are then mapped to target language sentences (if several derivations realize the same target sentence, their weights are summed) and reranked according to a language model (as was done in Huang et al. (2006)). This is the only part of the pipeline where we deviate from the theoretical log-linear model, and this is where we might make search errors. In principle, one could integrate the language model by intersection with the translation model (as the stateful MBOT model is closed under intersection with finite automata), but this is (currently) not computationally feasible due to the size of models.

Tuning Minimum error rate training (Och, 2003) is implemented using Z-MERT⁷ (Zaidan, 2009). A set of source sentences is (forest-)parsed and decoded; the translation forests are stored on disk. Then, in each iteration of Z-MERT, it suffices to extract k-best lists from the translation forests according to the current weight vector.

⁷<http://cs.jhu.edu/~ozaidan/zmert/>

5 WMT15 Experimental setup

We used the training data that was made available for the WMT15 shared translation task on English–German⁸. It consists of three parallel corpora (1.8M sentences of European parliament proceedings, 216K sentences of newswire text, and 2.3M sentences of web text after cleanup) and additional monolingual news data for language model training.

The English half of the parallel data was parsed using Egret⁹ which is a re-implementation of the Berkeley parser (Petrov et al., 2006). For the German parse, we used the BitPar parser (Schmid, 2004; Schmid, 2006). The BitPar German grammar is highly detailed, which makes the syntactic information contained in the parses extremely useful. Part-of-speech tags and category label are augmented by case, number and gender information, as can be seen in the German parse tree in Figure 1. We only kept the best parse for each sentence during training.

We then trained a 5-gram language model on monolingual data using KenLM¹⁰ (Heafield, 2011; Heafield et al., 2013). Word alignment was achieved using the `fast_align`¹¹ word aligner from `cdec` (Dyer et al., 2010). As usual, we discarded sentence pairs where one sentence was significantly longer than the other, as well as those that were too long or too short.

For tuning, we chose the WMT12 test set (3,003 sentences of newswire text), available as part of the development data for the WMT13 shared translation task. Since our system had limited coverage on this tuning set, we limited ourselves to the first a subset of sentences we could translate.

When translating the test set, our models used parse trees delivered by the Egret parser. After translation, recasing was done by examining the output syntax tree, using a simple heuristics looking for nouns and sentence boundaries as well as common abbreviations. Since coverage on the test set was also limited, we used a simple word-based fallback system whenever an untranslated state was encountered in a derivation tree.

⁸<http://www.statmt.org/wmt15/translation-task.html>

⁹<https://sites.google.com/site/zhangh1982/egret>

¹⁰<http://kheafield.com/code/kenlm/>

¹¹http://www.cdec-decoder.org/guide/fast_align.html

| BLEU | BLEU-cased | TER |
|------|------------|------|
| 15.3 | 14.4 | .777 |

Table 1: BLEU and TER scores of our system.

6 Results

We report the overall translation quality, as listed on <http://matrix.statmt.org/>, measured using BLEU (Papineni et al., 2002) and TER (Snover et al., 2006), in Table 1.

Results are significantly worse compared to last year’s system which used morphological enhancements such as compound splitting (Quernheim and Cap, 2014) and a phrase-based fallback system for sentences that the exact decoder could not handle. However, we should note that where the fallback system was not needed, we achieved a BLEU score of 16.7.

From a linguistic point of view, constructions that involve long-distance reordering and agreement are typically handled well. Figure 4 shows some example sentences from the WMT13 test set in comparison to a phrase-based baseline system.

On the other hand, our system frequently makes mistakes in lexical choice, and often uses rules that have been extracted from erroneous alignments. Sometimes, these mistakes cannot be alleviated by the language model due to data sparsity (no competing good candidate translation).

7 Conclusion and further work

We presented our submission to the WMT15 shared translation task based on a novel, promising “full syntax, no pruning” tree-to-tree approach to statistical machine translation, inspired by Huang et al. (2006). There are, however, still major drawbacks and open problems associated with our approach. Firstly, the coverage can still be significantly improved. In these experiments, our model was able to translate only 70% of the test sentences. To some extent, this number can be improved by providing more training data. Also, more rules can be extracted if we not only use the best parse for rule extraction, but multiple parse trees, or even switch to forest-based rule extraction (Mi and Huang, 2008). Finally, the size of the input parse forest plays a role. For instance, if we only supply the best parse to our model, translation will fail for approximately half of the input.

However, there are inherent coverage limits. Since our model is extremely strict, it will never

Verb missing:

- (M) wir haben zwei spezialisten für ihre stellungnahme *gebeten* .
 (“*we have two specialists for their statement asked .*”)
(P) wir haben zwei spezialisten für ihre stellungnahme .
 (“*we have two specialists for their statement .*”)
(R) wir haben die meinung von zwei fachärzten *eingeholt* .
(S) We *asked* two specialists for their opinion.

Plural noun with singular verb:

- (M) auch *das technische personal* hat mir sehr viel gebracht .
 (“*also the technical staff has me much brought .*”)
(P) auch *die technischen mitarbeiter* hat mir sehr viel gebracht .
 (“*also the technical co-workers has me much brought .*”)
(R) *das technische personal* hat mir ebenfalls viel gegeben .
(S) *The technical staff* has also brought me a lot.

No agreement between noun and adjective:

- (M) in diesem sinne werden die maßnahmen zum teil *das amerikanische demokratische system* untergraben .
 (“*in this sense will the measures to part (the american democratic system)_{NEUT} undermine .*”)
(P) in diesem sinne werden die maßnahmen teilweise , *die amerikanischen demokratische system* untergraben .
 (“*in this sense will the measures partially , the_{FEM} american democratic system_{NEUT} undermine .*”)
(R) in diesem sinne untergraben diese maßnahmen teilweise *das demokratische system der usa* .
(S) In this sense, the measures will partially undermine *the American democratic system*.

Long-distance reordering:

- (M) er zögert nicht , zu antworten , dass er einen antrag von einer unbekanntten person nie *akzeptieren würde* .
 (“*he hesitates not , to reply , that he a request from an unknown person never accept would .*”)
(P) er zögert nicht , sagen , dass er niemals *akzeptieren würde* einen antrag von einer unbekanntten person .
 (“*he hesitates not , say , that he never accept would a request from an unknown person .*”)
(R) gefragt antwortet er , dass er nie eine einladung von einem unbekanntten *annehmen würde* .
(S) He does not hesitate to reply that he *would* never *accept* a request from an unknown person.

Garbled output:

- (M) wie ich versprochen habe , ist meine tätigkeit teilweise reduziert worden .
 (“*as I promised have , has my activity partially reduced been .*”)
(P) wie ich ihnen zugesichert hatte , bestätigte , die meine aktivitäten wurden teilweise reduziert .
 (“*as I you assured had , confirmed , the my activities were partially reduced .*”)
(R) wie versprochen , habe ich meine aktivitäten teilweise zurückgefahren .
(S) As I promised, my activities have been partially reduced.

Figure 4: Examples from the test set where our ℓ MBOT system performed better, linguistically speaking; (M = ℓ MBOT system; P = phrase-based baseline system; R = reference translation; S = source sentence). Rough interlinear glosses are provided.

be able to translate sentences whose parse trees contain structures it has never seen before, since it has to match at least one input parse tree exactly. While we implemented a simple solution to handle unknown words, the issue with unknown structures is not so easy to solve without breaking the otherwise theoretically sound approach. Possibly, glue rules can help.

The second drawback is runtime. We were able to translate about 20 sentences per hour on one processor. Distributing the translation task on different machines, we were able to translate the WMT15 test set (3k sentences) in roughly three days. Given that the trend goes towards parallel programming, and considering the fact that our decoder is written in the rather slow language Python, we are confident that this is not a major problem. We were able to run the whole pipeline of training, tuning and evaluation on the WMT15

shared task data in less than one week. We are currently investigating whether A* k-best algorithms (Pauls and Klein, 2009; Pauls et al., 2010) can help to guide the translation process while maintaining optimality.

Thirdly, currently the language model is not integrated, but implemented as a separate reranking component. We are aware that an integrated language model might improve translation quality (see e.g. Chiang (2007) where 3–4 BLEU points are gained by LM integration). Some research on this topic already exists, e.g. (Rush and Collins, 2011) who use dual decomposition, and (Aziz et al., 2013) who replace intersection with an upper bound which is easier to compute. It might also be feasible to intersect the language model (represented by a regular string grammar) lazily.

References

- André Arnold and Max Dauchet. 1982. Morphismes et bimorphismes d'arbres. *Theoret. Comput. Sci.*, 20(1):33–93.
- Wilker Aziz, Marc Dymetman, and Sriram Venkatapathy. 2013. Investigations in exact inference for hierarchical translation. In *Proc. 8th WMT*, pages 472–483.
- Fabienne Braune, Nina Seemann, Daniel Quernheim, and Andreas Maletti. 2013. Shallow local multi-bottom-up tree transducers in statistical machine translation. In *Proc. 51th ACL*, pages 811–821.
- David Chiang. 2007. Hierarchical phrase-based translation. *Computat. Linguist.*, 33(2):201–228.
- Chris Dyer, Adam Lopez, Juri Ganitkevitch, Johnathan Weese, Ferhan Ture, Phil Blunsom, Hendra Setiawan, Vladimir Eidelman, and Philip Resnik. 2010. cdec: A decoder, alignment, and learning framework for finite-state and context-free translation models. In *Proc. 48th ACL*.
- Joost Engelfriet, Eric Lilin, and Andreas Maletti. 2009. Extended multi bottom-up tree transducers: Composition and decomposition. *Acta Inf.*, 46(8):561–590, October.
- Zoltán Fülöp and Heiko Vogler. 2009. Weighted tree automata and tree transducers. In Manfred Droste, Werner Kuich, and Heiko Vogler, editors, *Handbook of Weighted Automata*, EATCS Monographs on Theoret. Comput. Sci., chapter 9, pages 313–403. Springer.
- Kenneth Heafield, Ivan Pouzyrevsky, Jonathan H. Clark, and Philipp Koehn. 2013. Scalable modified Kneser-Ney language model estimation. In *Proc. 51st ACL*.
- Kenneth Heafield. 2011. KenLM: faster and smaller language model queries. In *Proc. 6th WMT*, pages 187–197.
- Liang Huang and David Chiang. 2005. Better k-best parsing. In *Proc. IWPT*, pages 53–64.
- Liang Huang, Kevin Knight, and Aravind Joshi. 2006. Statistical syntax-directed translation with extended domain of locality. In *Proc. 7th Conf. AMTA*, pages 66–73.
- Alexander Koller and Marco Kuhlmann. 2011. A generalized view on parsing and translation. In *Proc. IWPT*, pages 2–13.
- Eric Lilin. 1978. *Une généralisation des transducteurs d'états finis d'arbres: les S-transducteurs*. Thèse 3ème cycle, Université de Lille.
- Yang Liu, Yajuan Lü, and Qun Liu. 2009. Improving tree-to-tree translation with packed forests. In *Proc. 47th ACL*, pages 558–566.
- Andreas Maletti. 2010. Why synchronous tree substitution grammars? In *Proc. HLT-NAACL*, pages 876–884.
- Andreas Maletti. 2011. How to train your multi bottom-up tree transducer. In *Proc. 49th ACL*, pages 825–834.
- Andreas Maletti. 2012. Every sensible extended top-down tree transducer is a multi bottom-up tree transducer. In *Proc. HLT-NAACL*, pages 263–273.
- Haitao Mi and Liang Huang. 2008. Forest-based translation rule extraction. In *Proc. EMNLP*, pages 206–214.
- Haitao Mi, Liang Huang, and Qun Liu. 2008. Forest-based translation. In *Proc. 46th ACL*, pages 192–199. ACL.
- Franz Josef Och. 2003. Minimum error rate training in statistical machine translation. In *Proc. 41st ACL*, pages 160–167.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei jing Zhu. 2002. BLEU: a method for automatic evaluation of machine translation. In *Proc. 40th ACL*, pages 311–318.
- Adam Pauls and Dan Klein. 2009. K-best A* parsing. In *Proc. 47th ACL*, pages 958–966.
- Adam Pauls, Dan Klein, and Chris Quirk. 2010. Top-down k-best A* parsing. In *Proc. 48th ACL*, pages 200–204.
- Slav Petrov, Leon Barrett, Romain Thibaux, and Dan Klein. 2006. Learning accurate, compact, and interpretable tree annotation. In *Proc. COLING-ACL*, pages 433–440.
- Daniel Quernheim and Fabienne Cap. 2014. Large-scale exact decoding: The ims-ttt submission to wmt14. In *Proc. 9th WMT*, pages 163–170.
- Alexander M. Rush and Michael Collins. 2011. Exact decoding of syntactic translation models through lagrangian relaxation. In *Proc. 49th ACL*, pages 72–82.
- Helmut Schmid. 2004. Efficient parsing of highly ambiguous context-free grammars with bit vectors. In *Proc. 20th COLING*, pages 162–168.
- Helmut Schmid. 2006. Trace prediction and recovery with unlexicalized PCFGs and slash features. In *Proc. 44th ACL*.
- Matthew Snover, Bonnie Dorr, Richard Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *Proc. AMTA*.
- Jun Sun, Min Zhang, and Chew Lim Tan. 2009. A non-contiguous tree sequence alignment-based model for statistical machine translation. In *Proc. 47th ACL*, pages 914–922.

Omar F. Zaidan. 2009. Z-MERT: A fully configurable open source tool for minimum error rate training of machine translation systems. *The Prague Bulletin of Mathematical Linguistics*, 91:79–88.

Min Zhang, Hongfei Jiang, Aiti Aw, Haizhou Li, Chew Lim Tan, and Sheng Li. 2008a. A tree sequence alignment-based tree-to-tree translation model. In *Proc. 46th ACL*, pages 559–567.

Min Zhang, Hongfei Jiang, Haizhou Li, Aiti Aw, and Sheng Li. 2008b. Grammar comparison study for translational equivalence modeling and statistical machine translation. In *Proc. 22nd COLING*, pages 1097–1104.