

Dependency Treelet Translation: The convergence of statistical and example-based machine- translation?

Christopher Quirk and Arul Menezes

*Microsoft Research
One Microsoft Way
Redmond, WA 98052*

Abstract.

We describe a novel approach to machine translation that combines the strengths of the two leading corpus-based approaches: Phrasal SMT and EBMT. We use a syntactically informed decoder and reordering model based on the source dependency tree, in combination with conventional SMT models to incorporate the power of phrasal SMT with the linguistic generality available in a parser. We show that this approach significantly outperforms a leading string-based Phrasal SMT decoder and an EBMT system. We present results from two radically different language pairs, and investigate the sensitivity of this approach to parse quality by using two distinct parsers and oracle experiments. We also validate our automated BLEU scores with a small human evaluation.

1. Introduction

Data-Driven Machine Translation approaches, such as Example-Based Machine Translation and Statistical Machine Translation, have revolutionized the field of Machine Translation. Where once the only tractable approach toward MT was writing rule-based transfer systems, we are now seeing the emergence of large scale translation systems based on data-driven techniques. These approaches are founded on similar datasets but very different principles. However, we believe that that latest generation of DDMT systems demonstrate an increasing convergence of EBMT and SMT. After a brief survey of recent developments in EMBT and SMT, we describe a new syntax-based SMT approach that draws successfully from both traditions to produce high-quality translations.

Although it is difficult to pinpoint the exact definition of EBMT (for a more detailed discussion, the reader is referred to Carl (2006) and Wu (2006)), there are several fundamental principles that are generally accepted by the community as being characteristic of EBMT. First, EBMT is a data-driven approach: translation information is learned primarily from parallel corpora. Second, EBMT relies heavily on the concept of translation by analogy. When presented with a new sentence to translate, an EBMT system attempts to reuse translation information from its parallel corpus—preferably reusing information in segments that are as large as possible. In this sense, EBMT is

a generalization of Translation Memory: it goes beyond reuse of translations at the sentence level, attempting more aggressive reuse at the phrase, word, and perhaps morpheme or even character level. Lepage and Denoual (2006) is an example of a particularly pure EBMT system, relying on nothing except a pure analogic assembly mechanism and a parallel corpus. Translation examples are not limited to words or even contiguous phrases; often the most useful units of reuse are discontinuous in nature.

A common way for EBMT systems to exploit complex phenomena is through the use of linguistic analysis. As an example, consider the system of Menezes and Richardson (2003). The first step of training is to parse both source and target language sentences into a deep predicate-argument structure representation, which normalizes away many of the surface differences between languages. Next these deep structures are aligned; from this aligned structure, translation mappings can be automatically extracted. To translate a new sentence, it is first parsed into this deep representation. Then a target language deep representation is constructed by combining translation mappings learned from the parallel corpus. Finally the target language sentence is generated from the deep structure using a hand-written or machine-learned generation component. Other systems follow similar lines with some interesting variations. Kurohashi et al. (2005), for instance, obviates a target language generation component by employing a shallower analysis of both languages.

At the same time Statistical Machine Translation (SMT) began to develop using the same resources but different ideas. Parallel corpora also formed the foundation of SMT, though initial attempts at statistical translation models were less focused on the idea of reusing examples. Instead, effort was devoted to defining generative models that were sufficiently powerful to accommodate translational divergences while allowing tractable estimation. The touchstone of this early work in SMT is Brown et al. (1993), which defines a series of generative models each providing a distribution over the set of foreign language sentences and word alignments given a source language sentence. These models were to be used as channel models in a noisy channel decoder; n -gram language models (like those commonly used in speech recognition) could act as the target language model. However, these initial systems never achieved the speed or quality necessary to break into the mainstream translation market.

The recent activity in SMT has instead been driven by a byproduct of the generative models. Along with estimating parameters of channel models, these models also could be used to produce a word alignment. As the EBMT community had recognized for years, identifying word and phrase translations requires an a fine-grained correspondence between pieces of a sentence: an accurate word-alignment opened many possibilities for SMT. Seeing the potential power of larger translational units, approaches such as alignment templates (Och and Ney, 2004) and phrasal SMT (Koehn et al., 2003) took a

major step toward EBMT (perhaps an unintentional one) by extracting multi-word translations from a word-aligned parallel corpus and stringing them together to form a translation. However these new “phrasal” approaches were still grounded in statistical decision theory. Example phrase pairs with counts were not sufficient to form a full SMT system. The first systems still used something similar to the noisy channel approach to model translation, and tried to find heuristic search methods that approximated optimal decoding behaviour as much as possible. Latter systems generalized the decoding approach to form what are now called hybrid generative-discriminative models, using maximum entropy models (Och and Ney, 2002) or direct optimization of error rates (Och, 2003) to optimize functions.

With the above developments, one may easily argue that the convergence of EBMT and SMT had already begun. For instance, recent work, including Way and Gough (2005) and Groves and Way (2006), has noted and explored the similarities between phrasal SMT and Marker-based EBMT. Yet these phrasal SMT systems only started to exploit the information that had been used for years in EBMT systems. Even the very first EBMT systems found that effective reuse often requires non-contiguous phrases in either the source language or the target language. Other well-known phenomena (such as boundary friction (Somers, 2003)) also pose problems for these systems (though target language models may mitigate this somewhat). Nor does phrasal SMT have an explicit model to account for ordering differences between languages. While arbitrary reordering of words is allowed within memorized phrases, typically only a small amount of phrase reordering is allowed, often modeled in terms of string-level offsets. This reordering model is very limited in terms of linguistic generalizations. For instance, when translating English to Japanese, an ideal system would automatically learn large-scale typological differences: English SVO clauses generally become Japanese SOV clauses, English post-modifying prepositional phrases become Japanese pre-modifying postpositional phrases, etc. A phrasal SMT system may learn the internal reordering of specific common phrases, but it cannot generalize to unseen phrases that share the same linguistic structure.

A natural solution to many of these problems is the incorporation of syntax. Whether by incorporating linguistic parsers (as in Yamada and Knight (2002) or Yamada and Knight (2002)) or by simply including algorithms motivated by parsing (as in Wu (1997) or Chiang (2005)), syntax-based SMT allows natural incorporation of discontinuous phrases. Furthermore, syntax-based SMT also provides a natural means of modeling constituent reordering.

An early, elegant approach to syntax-based SMT is that of Inversion Transduction Grammars (Wu, 1997). No overt linguistic information is used; instead, it uses algorithms inspired by parsers. Translation is viewed as a process of parallel parsing of the source and target language via a synchronized grammar. To make this process computationally efficient, however, some se-

vere simplifying assumptions are made, such as using a single non-terminal label. This results in the model simply learning a very high level preference regarding how often nodes should switch order without any contextual information. Also these translation models are intrinsically word-based; phrasal combinations are not modeled directly, and results have not been competitive with the top phrasal SMT systems. Melamed (2004) generalizes this work, defining parsing algorithms over multitext grammars and demonstrating how these versatile tools can be used in various stages of a syntax-based SMT system.

In a similar vein, the Hiero approach (Chiang, 2005) also uses a decoding algorithm motivated by parsing, but incorporates a phrasal translations into translation. This is one of the first syntax-based SMT approaches to show significant improvements over a phrasal SMT baseline. There are several important lessons to be drawn from this work. Phrasal translations, effective decoding algorithms, and combinations of statistical models together produce a formidable backbone of an MT system.

Other systems do take advantage of actual linguistic information in translation. Yamada and Knight (2002) employ a parser in the target language to train probabilities on a set of operations that convert a English tree to a Japanese string. These channel models are then used in noisy-channel decoder to find the best English translation. Such an approach improves fluency slightly (Charniak et al., 2003), but does not significantly impact overall translation quality. This may be because the parser is applied to MT output, which is notoriously unlike native language, and no additional insight is gained via source language analysis. In a similar vein, Graehl and Knight (2004) proposes an framework for tree transduction that allows efficient training of transducers using the Expectation-Maximization (Dempster et al., 1977) algorithm. This promising approach of training channel models has the potential to improve string-to-tree translation.

Parsers can be also used in the source language. For instance, Lin (2004) translates dependency trees using paths. This is the first attempt to incorporate large phrasal SMT-style memorized patterns together with a separate source dependency parser and SMT models. However the phrases are limited to linear paths in the tree, the only SMT model used is a maximum likelihood channel model and there is no ordering model. Reported BLEU scores are not yet at the level of leading phrasal SMT systems.

Other approaches take advantage of both source and target language parsers. Imamura et al. (2005) uses both Japanese and English parsers to help limit the computational complexity and make certain linguistic generalizations more evident. This approach succeeds in outperforming a phrasal SMT baseline in a Japanese to English translation task using only a small set of models.

2. Dependency Treelet Translation

In this paper we propose a novel dependency tree-based approach to phrasal SMT that uses tree-based ‘phrases’ and a tree-based ordering model in combination with conventional SMT models to produce translations significantly better than a leading string-based system. We believe this approach reinforces the convergence of EBMT and SMT: our translations are firmly grounded in the training data, yet the translation process is guided by a number of probabilistic models in a general log-linear framework.

Our system employs a source-language dependency parser, a target language word segmentation component, and an unsupervised word alignment component to learn treelet translations from a parallel sentence-aligned corpus. We begin by parsing the source text to obtain dependency trees and word-segmenting the target side, then applying an off-the-shelf word alignment component to the bitext.

The word alignments are used to project the source dependency parses onto the target sentences. From this aligned parallel dependency corpus we extract a treelet translation model incorporating source and target treelet pairs, where a *treelet* is defined to be an arbitrary connected subgraph of the dependency tree. We also train a variety of statistical models on this aligned dependency tree corpus, including a channel model and an order model.

To translate an input sentence, we parse the sentence, producing a dependency tree for that sentence. We then employ a decoder to find a combination and ordering of treelet translation pairs that cover the source tree and are optimal according to a set of models that are combined in a log-linear framework as in Och and Ney (2003).

This approach offers the following advantages over string-based SMT systems: Instead of limiting learned phrases to contiguous word sequences, we allow translation by all possible phrases that form connected subgraphs (treelets) in the source and target dependency trees. This is a powerful extension: the vast majority of surface-contiguous phrases are also treelets of the tree; in addition, we gain discontinuous phrases, including combinations such as verb-object, article-noun, adjective-noun etc. regardless of the number of intervening words.

Another major advantage is the ability to employ more powerful models for reordering source language constituents. These models can incorporate information from the source analysis. For example, we may model directly the probability that the translation of an object of a preposition in English should precede the corresponding postposition in Japanese, or the probability that a pre-modifying adjective in English translates into a post-modifier in French.

2.1. CORPUS ANALYSIS AND PREPARATION

The following sections describe the process by which a word-aligned parallel dependency tree corpus is constructed from sentence-aligned data. These aligned parallel dependency trees are used to train the translation system.

2.1.1. *Source analysis*

We only require a relatively shallow source language analysis. We assume that source language fragments can be part-of-speech tagged, and a simple dependency analysis can be produced. Arc labels are not required. The dependency analysis is viewed in one of two isomorphic ways. First, we can look at this analysis as head annotation: each word in the sentence has a unique parent, except for one word, which is the root of the sentence. This parent function forms directed acyclic graph. Secondly, we can view this analysis in a tree-like manner: each word is annotated with a list of its premodifying children and its postmodifying children.

2.1.2. *Word alignment*

We also require a word alignment of the parallel corpus. A word alignment can be represented as a binary relation \sim between the source words and the target words in each sentence. The only restriction we place on this relation is that it cannot be many-to-many. That is, if S and T are sets of source and target words such that $s \sim t$ for all $s \in S$ and $t \in T$, then either $|S| = 1$ or $|T| = 1$.

However we currently obtain word alignments with GIZA++ (Och and Ney, 2003), which limits the word alignments to involve at most one word in the target side. Therefore we follow the common practice of deriving many-to-many alignments by running the IBM models in both directions and combining the results heuristically. Our heuristics differ in that they constrain many-to-one alignments to be contiguous in the source dependency tree. We apply the following rules in order:

1. Accept all alignments from the intersection.
2. Accept all alignments that are unique on both sides (i.e. the only alignment in the union from the given source word is to the given target word, and vice versa).
3. Accept all alignments that are unique on one side (i.e. the only alignment in the union from the given source word is to that target word, but the target word has other non-unique alignments, or vice versa).
4. Accept those many-to-one alignments that are adjacent to existing alignments in the source dependency tree (i.e. accept an alignment (s_i, t_k) if

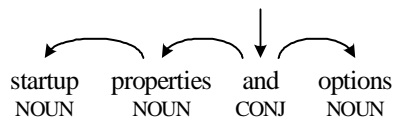


Figure 1. An example dependency tree.

we've already accepted an alignment (s_j, t_k) , and either the parent of s_i is s_j , or the parent of s_j is s_i .

5. Accept all one-to-many alignments.

The resulting alignment is a superset of the intersection and a subset of the union.

2.1.3. Dependency projection

Given a word aligned sentence pair and a source dependency tree, we use the alignment to project the source structure onto the target sentence. One-to-one alignments project directly to create a target tree isomorphic to the source. Many-to-one alignments project similarly; since the 'many' source nodes are connected in the tree, they act as if condensed into a single node. In the case of one-to-many alignments we project the source node to the rightmost¹ of the 'many' target words, and make the rest of the target words dependent on it.

Unaligned target words² are attached into the dependency structure as follows: assume there is an unaligned word t_j in position j . Let $i < j$ and $k > j$ be the target positions closest to j such that t_i depends on t_k or vice versa: attach t_j to the lower of t_i or t_k . If all the nodes to the left (or right) of position j are unaligned, attach t_j to the left-most (or right-most) word that is aligned. Algorithm 1 provides a pseudocode description of the process.

The target dependency tree created in this process may not read off in the same order as the target string, since our alignments do not enforce phrasal cohesion. For instance, consider the projection of the parse in Figure 1 using the word alignment in Figure 2a. Our algorithm produces the dependency tree in Figure 2b. If we read off the leaves in a left-to-right in-order traversal, we do not get the original input string: *de démarrage* appears in the wrong place.

A second reattachment pass corrects this situation. For each node in the wrong order, we reattach it to the lowest of its ancestors such that it is in the correct place relative to its siblings and parent. In Figure 2c, reattaching *démarrage* to *et* suffices to produce the correct order.

Algorithm 1 Tree projection algorithm

```

function ProjectDeptree(  $S$  : source nodes,  $T$  : target nodes,  $A$  : alignment)
   $a \leftarrow$  AlignmentFunction( $S, T, A$ )
   $h_t \leftarrow$  ProjectDeptreeBackbone(  $\varepsilon$ , root( $S$ ),  $\varepsilon$ ,  $a$ )
   $h_t \leftarrow$  AttachOthers( $h_t, T$ )
   $h_t \leftarrow$  Reattach( $t, h_t$ )
  return  $h_t$ 
function AlignmentFunction( $S, T, A$ )
   $a \leftarrow \emptyset$ 
  for all  $s \in S$  do
     $X \leftarrow \{t \in T \mid (s, t) \in A\}$ 
     $a(s) \leftarrow$  rightmost word in  $X$ ;  $\varepsilon$  if  $X = \emptyset$ 
  return  $a$ 
function ProjectDeptreeBackbone( $s_0, s_1, t_0, a$ )
   $h \leftarrow \emptyset$ 
   $t_1 \leftarrow a(s_1)$ 
  if  $t_0 \neq \varepsilon \wedge t_1 \neq \varepsilon$  then
     $h(t_1) \leftarrow t_0$ 
  for all  $s_2 \in \text{children}(s_1)$  do
     $h \leftarrow h \cup$  ProjectDeptreeBackbone( $s_1, s_2, t_1, a$ )
  return  $h$ 
function AttachOthers( $T, h, A$ )
  for all  $t \in T$  such that  $h(t) = \varepsilon$  do
    if  $t$  is aligned then
      Find  $s \in S, t' \in T$  such that  $(s, t), (s, t') \in A$  and  $h(t') \neq \varepsilon$ 
       $h(t) = t'$ 
    else
      Find closest aligned words to the left and right  $t_l, t_r$ 
       $h(t) \leftarrow t_l$  if  $t_l$  is further from the root than  $t_r$ ;  $h(t) \leftarrow t_r$  otherwise.
  return  $h$ 
function Reattach( $t, h$ )
   $Q \leftarrow \langle t \rangle$ ; a queue.
   $h' \leftarrow \emptyset$ ; the new parent relation without crossing dependencies
  while  $Q$  is not empty do
     $t_1 \leftarrow \text{Pop}(Q)$ 
     $t_0 \leftarrow h(t)$ 
    while  $(t_0, t_1)$  crosses some dependency already in  $h'$  do
       $t_0 \leftarrow h'(t_0)$ 
     $h'(t_1) \leftarrow t_0$ 
    Enqueue( $Q, \text{children}(t)$ )
  return  $h'$ 

```

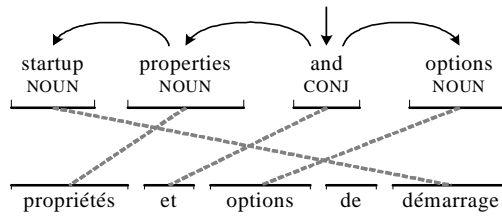


Figure 2a. Word alignment.

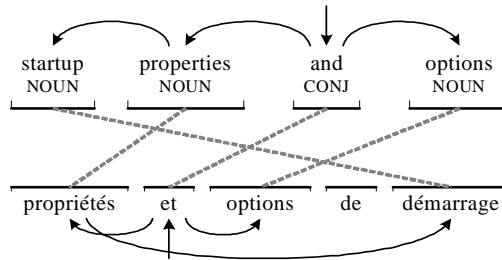


Figure 2b. Dependencies after initial projection step.

2.1.4. Extracting treelet translation pairs

From the aligned pairs of dependency trees we extract all pairs of aligned source and target treelets along with word-level alignment linkages, up to a configurable maximum size. We also keep treelet counts for maximum likelihood estimation.

2.2. MODELS

Generally we view translation as a global search problem: given an input dependency tree, a search component (or *decoder*) produces possible translation candidates, which are scored by a log-linear combination of feature functions. The highest scoring candidate is returned as the final translation.

A candidate in this search consists of:

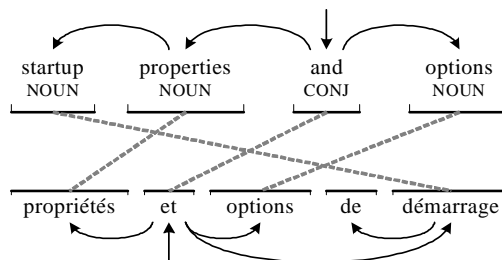


Figure 2c. Dependencies after reattachment.

- S : a source dependency tree,
- T : a target dependency tree,
- A : a word alignment between the source and target trees, and
- I : a set of treelet translation pairs that are a partitioning of S and T into treelets.

Put formally, then, we wish to find

$$\operatorname{argmax}_{T,A,I} \{\operatorname{SCORE}(S,T,A,I)\}$$

$\operatorname{SCORE}(S,T,A,I)$ is defined as a log-linear combination of the values of a set of feature functions F :

$$\operatorname{SCORE}(S,T,A,I) = \sum_{f \in F} w_f \cdot \log f(S,T,A,I)$$

Theoretically, these feature functions could be any real valued function of the candidate. Most, however, are simply the scores from probabilistic models. The following sections explore these models in more detail.

2.2.1. Order model

Phrasal SMT systems often use a model to score the ordering of a set of phrases. One approach is to penalize any deviation from monotone decoding; another is to estimate the probability that a source phrase in position i translates to a target phrase in position j (Koehn et al., 2003).

We attempt to improve on these approaches by incorporating syntactic information into the ordering process. Our model assigns a probability to the order of a target tree given an unordered target tree, a source tree, and a word alignment between the two. Since we assume that dependencies are projective and phrases cohere during translation (i.e., the translation of a source constituent is a target constituent that is a contiguous substring), the location of a constituent within a sentence is determined by the position of its root node relative to its parent and the other modifiers of that parent. One way to indicate this order is with head-relative positions: the closest premodifier of a node has position -1 , the next has position -2 , and so on; the closest postmodifier of a node has position $+1$, and so on. Figure 3 demonstrates an aligned dependency tree pair annotated with head-relative positions.

We can estimate a conditional probability distribution $\Pr(T|S,Q,A)$ over ordered target trees T given a source tree S , an unordered target tree Q (a tree with a head function, but no relative order between nodes), and a the word alignment A between S . Then, given a function ϕ from ordered trees to their

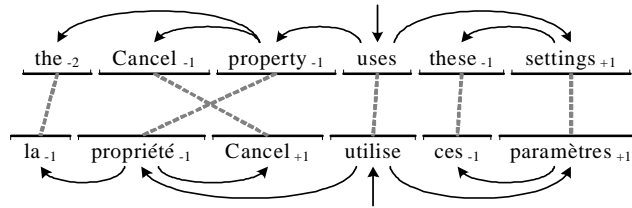


Figure 3. Aligned dependency tree pair annotated with head-relative positions.

unordered counterparts, we can use this probabilistic order model as another feature function in the log-linear combination:

$$f_{\text{Order}}(S, T, A, I) = \mathbf{Pr}(T|S, \varphi(T), A)$$

We model this distribution in terms of head relative orderings, and we make the strong independence assumption that each node is positioned independently:

$$\mathbf{Pr}(T|S, Q, A) \approx \prod_{t \in Q} \mathbf{Pr}(\text{pos}(t)|S, Q, A)$$

Furthermore, to make training tractable and to simplify decoding, we use only a very small set of features reflecting local information in the dependency tree to model this probability. In particular, these are the features used:

- The lexical items of the target node and its parent.
- The lexical items of the source nodes aligned to the target node and its parent.
- The part-of-speech (category, or “cat”) of the source nodes aligned to the node and its parent.
- The head-relative position of the source node aligned to the node.

One could also make a Markov assumption and condition on information from the siblings of the target node; we found that this complicated the model and had little impact in practice.

As an example, consider the children of *propriété* in Figure 3. The head-relative positions of its modifiers *la* and *Cancel* are -1 and +1, respectively. Then the model attempts to predict the target positions as described in Figure 4. Despite all the limitations of this model, it can capture important linguistic generalizations that are not available to purely phrase-based systems. For instance, the generalization that English premodifiers of Nouns generally become postmodifiers in French unless they are determiners is straightforward to learn from a decision tree over the given features; similarly, one can see that English prepositions generally become Japanese postpositions.

$$\begin{aligned}
\Pr(\text{pos}(m1) = -1 \mid & \text{lex}(m1) = \textit{la}, \text{lex}(h) = \textit{propriété}, \\
& \text{lex}(\text{src}(m1)) = \textit{the}, \text{lex}(\text{src}(h)) = \textit{property}, \\
& \text{cat}(\text{src}(m1)) = \text{DET}, \text{cat}(\text{src}(h)) = \text{NOUN}, \\
& \text{pos}(\text{src}(m1)) = -2) \cdot \\
\Pr(\text{pos}(m2) = +1 \mid & \text{lex}(m2) = \textit{Cancel}, \text{lex}(h) = \textit{propriété}, \\
& \text{lex}(\text{src}(m2)) = \textit{Cancel}, \text{lex}(\text{src}(h)) = \textit{property}, \\
& \text{cat}(\text{src}(m2)) = \text{NOUN}, \text{cat}(\text{src}(h)) = \text{NOUN}, \\
& \text{pos}(\text{src}(m2)) = -1)
\end{aligned}$$

Figure 4. Example order model decomposition

The training corpus in the form of word-aligned parallel dependency trees acts as a training set for supervised classifiers. From each target language node, we extract a single data point, where the target feature is the head-relative position in the target tree, and the other features are simply those available in the training set.

Such a training set could be used with many different machine learning methods. We have found decision trees to be surprisingly effective: they are both fast to train, relatively robust, and able to deal discover feature conjunctions, hence solve problems that are not linearly separable.

For a given test feature vector, we compute a probability distribution from the decision tree by first following the path to the leaf dictated by that feature vector, and then use the maximum likelihood estimate over the target feature counts at that leaf (Chickering, 2002).

2.2.2. Channel models

We incorporate two distinct channel models, a maximum likelihood estimate (MLE) model and a model computed using Model-1 word-to-word alignment probabilities as in Vogel et al. (2003). The MLE model effectively captures non-literal phrasal translations such as idioms, but suffers from data sparsity. The word-to-word model does not typically suffer from data sparsity, but prefers more literal translations.

Given a set of treelet translation pairs that cover a given input dependency tree and produce a target dependency tree, we model the probability of source given target as the product of the individual treelet translation probabilities: we assume a uniform probability distribution over the decompositions of a tree into treelets. We have four channel model feature functions, since a model can either predict target given source (a direct model) or source given target (an inverse model), and it can be estimated using maximum likelihood or via

lexical translation probabilities.

$$\begin{aligned}
 f_{\text{DirectMLE}}(S, T, A, I) &= \prod_{\langle \sigma, \tau \rangle \in I} \frac{c(\sigma, \tau)}{c(\sigma, *)} \\
 f_{\text{InverseMLE}}(S, T, A, I) &= \prod_{\langle \sigma, \tau \rangle \in I} \frac{c(\sigma, \tau)}{c(*, \tau)} \\
 f_{\text{DirectLex}}(S, T, A, I) &= \prod_{\langle \sigma, \tau \rangle \in I} \prod_{t \in \tau} \sum_{s \in \sigma} p(s|t) \\
 f_{\text{InverseLex}}(S, T, A, I) &= \prod_{\langle \sigma, \tau \rangle \in I} \prod_{s \in \sigma} \sum_{t \in \tau} p(t|s)
 \end{aligned}$$

The lexical translation probabilities $p(s|t)$ and $p(t|s)$ used here are from Model 1 of Brown et al. (1993).

2.2.3. Target language model

One crucial component of modern statistical translation systems is a target language model. These models appear simple — model the probability of a target string with a Markov assumption, i.e. find the probability of each word in the context of the previous n words — but have a major impact on the fluency, grammaticality, and even accuracy of translation. As our dependency trees are ordered, an in-order walk suffices to enumerate the target words in order, so adding an n -gram target language model as another feature is quite easy:

$$f_{\text{Target}}(S, T, A, I) = \prod_{i=1}^{|T|} \Pr(t_i | t_{i-n}^{i-1})$$

We estimate the n -gram probabilities using modified Kneser-Ney smoothing (Goodman, 2001).

2.2.4. Miscellaneous feature functions

The log-linear framework allows us to incorporate other feature functions as ‘models’ in the translation process. For instance, using fewer, larger treelet translation pairs often provides better translations, since they capture more context and allow fewer possibilities for search and model error. Therefore we add a feature function that counts the number of phrases used. We also add a feature that counts the number of target words; this acts as an insertion/deletion bonus/penalty, and counteracts the preference of the target language model for shorter sentences.

$$\begin{aligned}
 f_{\text{PhraseCount}}(S, T, A, I) &= e^{|I|} \\
 f_{\text{WordCount}}(S, T, A, I) &= e^{|T|}
 \end{aligned}$$



Figure 5. Example input dependency tree

2.3. DECODING

The challenge of tree-based decoding is that the traditional left-to-right decoding approach of string-based systems is inapplicable. Additional challenges are posed by the need to handle treelets—perhaps discontinuous or overlapping—and a combinatorially explosive ordering space.

Our decoding approach is influenced by ITG (Wu, 1997) with several important extensions. First, we employ treelet translation pairs instead of single word translations. Second, instead of modeling rearrangements as either preserving source order or swapping source order, we allow the dependents of a node to be ordered in any arbitrary manner and use the order model to estimate probabilities. Finally, we use a log-linear framework for model combination that allows any amount of other information to be modeled.

We will initially approach the decoding problem as a bottom up, exhaustive search. We define the set of all possible treelet translation pairs of the subtree rooted at each input node in the following manner: A treelet translation pair x is said to *match* the input dependency tree S iff there is some connected subgraph S' that is identical to the source side of x . We say that x *covers* all the nodes in S' and is *rooted* at source node s , where s is the root of matched subgraph S' .

Consider in turn each treelet translation pair x rooted at s . The treelet pair x may cover only a portion of the input subtree rooted at s . Find all descendants s' of s that are not covered by x , but whose parent s'' is covered by x . At each such node s'' look at all interleavings of the children of s'' specified by x , if any, with each translation t' from the candidate translation list of each child s' , which is computed in the same way on-demand and memoized. Each such interleaving is scored using the models previously described and added to the candidate translation list for that input node. The resultant translation is the best scoring candidate for the root input node.

As an example, see the example dependency tree in Figure 5 and treelet translation pair in Figure 6. This treelet translation pair covers all the nodes in 5 except the subtrees rooted at *software* and *is*. We first compute (and cache) the candidate translation lists for the subtrees rooted at *software* and *is*, then construct full translation candidates by attaching those subtree translations to *installés* in all possible ways. The order of *sur* relative to *installés* is fixed; it remains to place the translated subtrees for the *software* and *is*. Note that if c is the count of children specified in the mapping and r is the count of subtrees

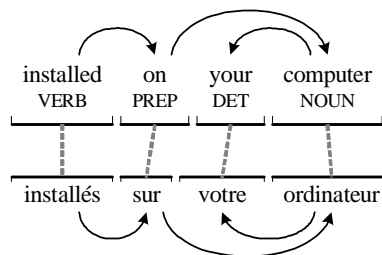


Figure 6. Example treelet translation pair

translated via recursive calls, then there are $(c + r + 1)! / (c + 1)!$ orderings. Thus $(1 + 2 + 1)! / (1 + 1)! = 12$ candidate translations are produced for each combination of translations of *the software* and *is*.

2.3.1. Dynamic Programming

Converting this exhaustive search to dynamic programming relies on several key observations. First we note that the candidate scores are generally multiplicative. The channel model probability of a candidate constructed from a treelet and existing candidates for subtrees not covered by that treelet is simply the product of the treelet and candidate channel model probabilities; a similar situation exists for the word count and phrase count features. For the order model, we must also score the head relative positions of the treelet as well as the root elements of the existing candidates, and for the target language model, we must multiply in the probabilities of the words at the boundary of each candidate. Secondly, we note that these additional probabilities only depend on a very small amount of information in the candidate: the order model probability requires only on the lexical item at the root of the tree, and a n -gram target language model is requires only the first and last $(n - 1)$ words in each subtree. Therefore, for any given source subtree, we only need to keep the best scoring translation candidate for each combination of (head, leading $(n - 1)$ -gram, $(n - 1)$ -gram bigram). In the worst case, however, this still requires keeping $O(V^5)$ candidates per input subtree, where V is the target language vocabulary size.

Although dynamic programming does limit the search space, it alone is not sufficient to produce a real-time translation system. Therefore we explore the following optimizations that have proven effective in experimental settings, but do not preserve optimality.

2.3.2. Beam search

Instead of keeping the full list of translation candidates for a given input node, we keep a top-scoring subset of the candidates. While the decoder is no longer

guaranteed to find the optimal translation, in practice the quality impact is minimal with a beam size of approximately 10.

2.3.2.1. Variable-sized n -best lists A further speedup can be obtained by noting that the number of translations using a given treelet pair is exponential in the number of subtrees of the input not covered by that pair. To limit this explosion we vary the beam size in inverse proportion to the number of subtrees uncovered by the current treelet. This has the intuitive appeal of allowing a more thorough exploration of large treelet translation pairs (that are likely to result in better translations) than of smaller, less promising pairs.

2.3.3. Pruning treelet translation pairs

Channel model scores and treelet size are powerful predictors of translation quality. Heuristically pruning low scoring treelet translation pairs before the search starts allows the decoder to focus on combinations and orderings of high quality treelet pairs.

- Only keep those treelet translation pairs with an MLE probability above a threshold t . This has the greatest impact on small, common treelets (such as translations of the English word “the”) that have many possible translations. The vast majority of these translations, however, are due to a single sentence pair in which the word alignment was suboptimal.
- Given a set of treelet translation pairs with identical sources, keep those with an MLE probability within a ratio r of the best pair.
- At each input node, keep only the top k treelet translation pairs rooted at that node, as ranked first by size, then by MLE channel model score, then by Model 1 score. This final pruning heuristic has the greatest impact on translation speed.

2.3.4. Greedy ordering

The complexity of the ordering step at each node grows with the factorial of the number of children to be ordered. This can be tamed by noting that given a fixed pre- and post-modifier count, our order model is capable of evaluating a single ordering decision independently from other ordering decisions.

One version of the decoder takes advantage of this to severely limit the number of ordering possibilities considered. Instead of considering all interleavings, it considers each potential modifier position in turn, greedily picking the most probable child for that slot, moving on to the next slot, picking the most probable among the remaining children for that slot and so on.

The complexity of greedy ordering is linear, but at the cost of a noticeable drop in BLEU score (see results in the experimentation section). Under de-

Table I. Data characteristics

		English	French	English	Japanese
Training	Sentences	500,000		500,000	
	Words	6,598,914	7,234,153	7,909,198	9,379,240
	Vocabulary	72,440	80,758	66,731	68,048
	Singletons	38,037	39,496	50,381	52,911
Test	Sentences	10,000		10,000	
	Words	133,402	153,701	175,665	211,139

fault settings our system tries to decode a sentence with exhaustive ordering until a specified timeout, at which point it falls back to greedy ordering.

3. Technical data experiments

We evaluated the translation quality of the system using the BLEU metric (Papineni et al., 2002) under a variety of configurations. We compared against two radically different types of systems to demonstrate the competitiveness of this approach:

- Pharaoh: A leading phrasal SMT decoder (Koehn et al., 2003).
- The MSR-MT system described in Section 1, an EBMT/hybrid MT system.

3.1. LANGUAGE PAIRS

We ran experiments in translating English to French and English to Japanese. The latter was chosen deliberately to highlight the challenges facing string-based MT approaches in language pairs with significant word-order differences.

Word order in Japanese is fundamentally very different from English. English is generally SVO (subject first, then verb, then object), where Japanese is SOV with a strong bias for head-final structures. Several other differences include:

- Word order is more flexible, since verbal arguments are generally indicated by postpositions, e.g. a direct object is indicated by the postposition ‘o’, a subject by ‘ga’.

- Many English phrases that are realized as post-modifiers (such as relative clauses and prepositional phrases) are translated as Japanese pre-modifiers; demonstratives and adjectives remain pre-modifiers.
- Verbal and adjectival morphology in Japanese is relatively complex: information contained in English pre-modifying modals and auxiliaries is often represented as verbal morphology.
- Japanese nouns and noun phrases are not marked for definiteness or number.

The word-aligned sentence pair in Figure 5 demonstrates many of these phenomena.

3.2. DATA

We used a corpus of Microsoft technical data (e.g., support articles, product documentation) containing over 1 million sentence pairs for each language-pair. We excluded sentences containing XML or HTML tags and for each language pair randomly selected training data sets ranging from 1,000 to 500,000 sentence pairs as well as 10,000 sentences for development testing and parameter tuning, 250 sentences for lambda training and 10,000 sentences for testing. Table I presents basic characteristics of these corpora.

3.3. TRAINING

We parsed the source (English) side of the corpus using two different parsers: NLPWIN, a broad-coverage rule-based parser developed at Microsoft Research able to produce syntactic analyses at varying levels of depth (Heidorn, 2000), and a Treebank parser (Bikel, 2004). For the purposes of these experiments we used a dependency tree output with part-of-speech tags and unstemmed, case-normalized surface words.

For word alignment, we used GIZA++ (Och and Ney, 2003), following a standard training regimen of five iterations of Model 1, five iterations of the HMM Model, and five iterations of Model 4, in both directions. Treelets were extracted from this word aligned parallel corpus; Table II presents some statistics about those corpora. We note that fewer treelet translation pairs were extracted from the English-Japanese parallel corpus; this is presumably due to the difficulty of aligning that language pair.

Target language models were trained using only the French and Japanese sides, respectively, of the parallel corpus; additional monolingual data may improve its performance. Finally we trained lambdas via Maximum BLEU (Och, 2003) on 250 held-out sentences with a single reference translation, and tuned the decoder optimization parameters (beam size, cutoffs, etc.) on the development test set.

Table II. Treelet statistics

	English→French, 300k	English→Japanese, 500k
Total treelet translation pairs	42,201,316	35,420,445
Distinct treelet translation pairs	30,188,949	20,600,885
Distinct source treelets	12,659,291	9,554,323

Table III. System Comparisons

	English→French, 100k		English→Japanese, 500k	
	BLEU	Sents/min	BLEU	Sents/min
Pharaoh monotone	37.06	4286	25.06	1600
Pharaoh	38.83	162	30.58	82
MSR-MT	35.26	453	-	-
Treelet	40.66	10.1	33.18	21

3.3.1. *Pharaoh*

The same GIZA++ alignments as above were used in the Pharaoh decoder. We used the heuristic combination described in (Och and Ney, 2003) and extracted phrasal translation pairs from this combined alignment as described in (Koehn et al., 2003). Except for the order model (Pharaoh uses a penalty on the deviance from monotone), the same models were used: MLE channel model, Model 1 channel model, target language model, phrase count, and word count. Lambdas were trained in the same manner (Och, 2003).

3.3.2. *MSR-MT*

MSR-MT used its own word alignment approach as described in (Menezes and Richardson, 2003) on the same training data. MSR-MT does not use lambdas or a target language model.

3.4. RESULTS

We present BLEU scores on an unseen 10,000 sentence test set using a single reference translation for each sentence. Speed numbers are the end-to-end translation speed in sentences per minute. Unless otherwise specified all results are based on a phrase/treelet size of 4 and a training set size of 100,000 sentences for English to French and 500,000 sentences for English to Japanese. Unless otherwise noted all the differences between systems are statistically significant at $P < 0.01$.

Table IV. BLEU scores on training data subsets, phrase/treelet size 4

	1k	3k	10k	30k	100k	300k	500K
<i>English→French</i>							
Pharaoh	17.20	22.51	27.70	33.73	38.83	42.75	-
Treelet	18.70	25.39	30.96	35.81	40.66	44.32	-
<i>English→Japanese</i>							
Pharaoh	14.85	15.99	18.18	21.89	23.01	26.67	30.58
Treelet	13.90	15.39	18.94	23.99	25.68	29.97	33.18

Table V. Effect of maximum treelet/phrase size measured in BLEU scores

Max size	English→French		English→Japanese			
	100K		100K		500K	
	Treelet	Pharaoh	Treelet	Pharaoh	Treelet	Pharaoh
1	37.50	23.18	22.36	12.75	26.95	17.72
2	39.84	32.07	24.53	18.63	31.33	24.30
3	40.36	37.09	25.44	21.37	32.58	28.15
(default) 4	40.66	38.83	25.68	23.01	33.18	30.58
5	40.71	39.41	25.87	23.82	-	-
6	40.74	39.72	25.92	24.43	-	-

Comparative results are presented in Table III. Pharaoh monotone refers to Pharaoh with phrase reordering disabled.

Table IV compares the systems at different training corpus sizes. All differences are statistically significant at $P < 0.01$ except for English→Japanese at training set sizes less than 30K. Note that in English→French, where word order differences are mainly local, the gap between the systems narrows slightly with larger corpus sizes, however in English→Japanese, with global ordering differences, the treelet system’s margin over Pharaoh (initially negative) actually increases with increasing corpus size.

Table V compares Pharaoh and the Treelet system at different phrase sizes. The wide gap at smaller phrase sizes is particularly striking. It appears that while Pharaoh depends heavily on long phrases to encapsulate reordering, our dependency tree-based ordering model enables credible performance even with short phrases/treelets. Our treelet system with two-word treelets outperforms Pharaoh with six-word phrases.

Table VI presents some statistics on the number of treelets available during decoding time as well as the average number of treelets used in an individ-

Table VI. Treelet statistics at decoding time

	English→French 300K	English→Japanese 500K
Average source treelets matched	61.84	57.24
Average treelet translation pairs available	114.84	130.20
Average treelet translation pairs used	8.44	9.96

Table VII. Effect of ordering strategy, measured by BLEU score

	English→French 100K		English→Japanese 500K	
	BLEU	Sents/min	BLEU	Sents/min
<i>Monotone</i>				
Pharaoh	37.06	4286	25.06	1600
Treelet: no order model	35.35	39.7	26.43	67
<i>Non-monotone</i>				
Pharaoh	38.83	162	30.58	82
Treelet: greedy ordering	38.85	13.1	31.99	43
Treelet: exhaustive	40.66	10.1	33.18	21

ual sentence. The number of treelets matched per sentence is slightly less in English-Japanese vs. English-French; this is probably due to the smaller number of treelets extracted overall. However, we note that the number of treelet translation pairs is slightly higher. Although this could be due to a greater amount of ambiguity in English-Japanese translation, we think a more likely explanation is that the lower word alignment quality leads to less consistent mappings, and therefore more possible translation pairs. Finally, we see that on average more translation pairs are used to translate a single sentence; this may partly account for the lower BLEU scores in English-Japanese.

Table VII compares different ordering strategies. In contrast to results reported for English-Chinese (Vogel et al., 2003), monotone decoding severely degrades the performance of both systems in English→Japanese. We presume that this is due to the broad differences in word order between the two languages. In English→French the degradation is less marked.

Table VIII shows the translation results are not dependent on one particular parser, though a parser trained on a different domain (here, the Treebank) is at a disadvantage.

Table IX shows the impact of using parses beyond the 1-best. The first experiment is quite simple: we translate each parse separately, and keep the

Table VIII. Using different parsers (English→Japanese, 100K, size 4)

	BLEU
Pharaoh	23.01
NLPWIN parser: top parse only	25.68
Bikel parser: top parse only	24.15

Table IX. Using k -best parses (English→Japanese, 500k, size 4)

	BLEU
Pharaoh	30.58
Single NLPWIN parse	33.18
Top 100 NLPWIN parses	34.13
Oracle selection (top 100 NLPWIN parses)	36.91

1-best translation from each. Finally, we pick the highest scoring translation from each of these with no additional information about the parse. Even without any features of the parse itself, this approach boosts the BLEU score by a significant amount. This highlights the problems with the pipelined approach of using the single best parse, and suggests that decoding over a packed forest of parses may produce significantly better results.

The last line in Table IX is an oracle experiment to demonstrate the potential improvements from better parse selection. In this experiment, we translate each parse separately and again keep only the top scoring translation for each parse. However, instead of picking among these translations based on score, we consult the reference translation and pick the candidate most like the reference; the improvement in quality is quite impressive. Better parse ranking mechanisms may therefore have a positive impact on translation quality. Unfortunately we do not have source language sentences with both gold-standard dependency annotations and target language reference translations, so it is difficult to identify the correlation between monolingual parse accuracy and efficacy in translation.

Returning to the 1-best parse, in Table X we see a translation oracle experiment that demonstrates the impact of model error. The oracle picks the translation most like the reference translation from among the top n translations produced by the treelet system. Better models have the potential for major quality improvements, though Och et al. (2004) suggests achieving this gain is difficult.

Table X. Translation oracle (English→Japanese, 500k, size 4)

Translations	BLEU
1	33.18
4	35.30
16	37.38
64	38.56
256	38.70

Table XI. Human evaluation of 100 sentences (English→Japanese, 500k, size 4)

	Rater 1	Rater 2
Treelet preferred	30	50
Neither preferred	59	34
Pharaoh preferred	11	16

Finally, we expect the gains due to better parse selection may be somewhat orthogonal to the gains due to better channel model selection. Often when a parse is incorrect, it precludes selection of the correct translation. In the future we plan to explore the interactions between k -best parsing and n -best translations in more detail. Discriminative reranking including features from the source parse tree are particularly promising.

3.5. HUMAN EVALUATION

Two human raters were presented (in random order) both Pharaoh and Treelet translations of 100 sentences between 10 and 25 words and corresponding source and reference translations. They were asked to pick the more accurate translation. Table XI shows that for most of the sentences, humans prefer the Treelet translations, which is consistent with the BLEU scores above.

4. Conclusions and future work

We presented a novel approach to syntactically-informed statistical machine translation leveraging a parsed dependency tree representation of the source language via a tree-based ordering model and a syntactically informed decoder. We showed that it outperforms a leading phrasal SMT decoder in BLEU and human quality judgments. We also showed that it out-performed our own logical form-based EBMT/hybrid MT system.

SOURCE	in Visual Studio .NET , create a new Managed C++ application called determineOS .
REFERENCE	Visual Studio .NET で determineOS という名前で新しい Managed C アプリケーションを作成します。
PHARAOH	Visual Studio .NET では、新しい Managed C アプリケーション determineOS と呼ばれます。
TREELET	Visual Studio .NET で determineOS と呼ばれる新しい Managed C アプリケーションを作成します。
ANALYSIS	<i>Pharaoh drops main verb "create"; "call" becomes main verb; "determineOS" incorrectly compounded onto "application"</i>
<hr/>	
SOURCE	in the Named box , type scsi1hlp.vxd , and then click Find Now .
REFERENCE	[ファイルまたはフォルダの名前]ボックスに scsi1hlp.vxd と入力し、[検索 開始] ボタンをクリックします。
PHARAOH	[名前]ボックスに入力し、[検索 開始]をクリックします。 scsi1hlp.vxd
TREELET	[名前]ボックスで、 scsi1hlp.vxd、入力し[検索 開始]をクリックします。
ANALYSIS	<i>Pharaoh incorrectly moves "scsi1hlp.vxd" end of sentence. Treelet translation is not perfect, however: should have placed accusative marker on "scsi1hlp.vxd" and used "に" instead of "で".</i>
<hr/>	
SOURCE	you may be able to recover some disk space by quitting unneeded programs .
REFERENCE	不要なプログラムを終了して、ディスクの空き領域を回復します。
PHARAOH	一部のディスク領域を回復できる場合があります不要なプログラムを終了します。
TREELET	不要なプログラムを終了して一部のディスクスペースを回復できる場合があります。
ANALYSIS	<i>Pharaoh directly concatenates clauses with no conjunction.</i>
<hr/>	
SOURCE	when an ALTER TABLE statement is executed , the ROWCOUNT value of the session is taken into account .
REFERENCE	ALTER TABLE ステートメントが実行されるときに、そのセッションの ROWCOUNT の値が考慮されます。
PHARAOH	ALTER TABLE ステートメントを実行すると、セッションが ROWCOUNT の値にします。
TREELET	ALTER TABLE ステートメントが実行されると、ROWCOUNT の値を、セッションが考慮します。
ANALYSIS	<i>Pharaoh loses main verb "taken into account". Again Treelet translation is not perfect: less fluent word order, and the transitive verb has an inanimate subject, which is strongly dispreferred by Japanese speakers.</i>

Figure 7. Sentences where the Treelet translation was preferred over the Pharaoh translation.

SOURCE	type a new name for the Riched32.dll file (for example , Riched32.old) , and then press ENTER .
REFERENCE	Riched32.dll の新しい名前 (たとえば、Riched32.old) を入力して、Enter キーを押します。
PHARAOH	Riched32.dll ファイル (たとえば、Riched32.old の新しい名前を入力し、Enter キーを押します。
TREELET	Riched32.dll ファイル (たとえば、Riched32.old)、の新しい名前を入力とし、Enter キーを押します。
ANALYSIS	<i>Treelet incorrectly translates "type" as "入力とし" ("consider as input") rather than "入力し" ("input")</i>
<hr/>	
SOURCE	click to select the Use the Same Proxy Server for All Protocols check box .
REFERENCE	[すべてのプロトコルに同じプロキシサーバーを使用する] チェックボックスをクリックしてオンにします。
PHARAOH	同じ Proxy Server のすべてのプロトコルを使用する] チェックボックスをオンにします。
TREELET	クリック、チェックボックスをすべてのプロトコルに同じプロキシサーバーを使用して選択します。
ANALYSIS	<i>Treelet misparses input sentence: "Use the Same Proxy Server for All Protocols" is not identified as a user interface term modifying "check box". Main verb is therefore incorrect, amongst other errors.</i>

Figure 8. Sentences where the Treelet translation was preferred over the Pharaoh translation.

Even in the absence of a parse quality metric, we found that employing multiple parses could improve translation quality. Adding a parse probability may help further the gains from these additional possible analyses.

The syntactic information used in these models is still rather shallow. Order modeling may benefit from additional information such as semantic roles or morphological features. Furthermore, different model structures, machine learning techniques, and target feature representations all have the potential for significant improvements.

Notes

¹ If the target language is Japanese, leftmost may be more appropriate.

² Source unaligned nodes do not present a problem, with the exception that if the root is unaligned, the projection process produces a forest of target trees anchored by a dummy root.

Acknowledgements

We would like to thank Colin Cherry for his contributions and insight, Bob Moore for many enlightening technical discussions, and Chris Brockett, Hisami Suzuki, and Takako Aikawa for crucial insights and discussions of Japanese linguistic phenomena.

References

- Bikel, D. M.: 2004, 'Intricacies of Collins Parsing Model'. *Computational Linguistics* **30**(4), 479–511.
- Brown, P. F., S. A. D. Pietra, V. J. D. Pietra, and R. L. Mercer: 1993, 'The Mathematics of Statistical Machine Translation: Parameter Estimation'. *Computational Linguistics* **19**(2), 263–311.
- Carl, M.: 2006, 'A System-Theoretic View on EBMT'. *Machine Translation this volume*.
- Charniak, E., K. Knight, and K. Yamada: 2003, 'Syntax-based Language Models for Statistical Machine Translation'. In: *MT Summit*. New Orleans, Louisiana, USA, pp. 40–46.
- Chiang, D.: 2005, 'A hierarchical phrase-based model for statistical machine translation'. In: *Proceedings of the Annual Meeting of the Association for Computational Linguistics*. Ann Arbor, Michigan, USA, pp. 263–270.
- Chickering, D. M.: 2002, 'The WinMine Toolkit'. Technical Report MSR-TR-2002-103, Microsoft Research.
- Dempster, A., N. Laird, and D. Rubin: 1977, 'Maximum likelihood from incomplete data via the EM algorithm'. *Journal of the Royal Statistical Society* **39**(1), 1–38.
- Goodman, J.: 2001, 'A Bit of Progress in Language Modeling'. Technical Report MSR-TR-2001-72, Microsoft Research.
- Graehl, J. and K. Knight: 2004, 'Training Tree Transducers'. In: *Proceedings of HLT/NAACL*. Boston, Massachusetts, USA, pp. 105–112.
- Groves, D. and A. Way: 2006, 'Hybrid Data-Driven Models of Machine Translation'. *Machine Translation this volume*.
- Heidorn, G.: 2000, 'Intelligent writing assistance'. In: R. Dale, H. Moisl, and H. Somers (eds.): *Handbook of Natural Language Processing*. Marcel Dekker, pp. 181–208.
- Imamura, K., H. Okuma, and E. Sumita: 2005, 'Practical Approach to Syntax-based Statistical Machine Translation'. In: *MT Summit*. Phuket, Thailand, pp. 267–274.
- Koehn, P., F. J. Och, and D. Marcu: 2003, 'Statistical Phrase-Based Translation'. In: *Proceedings of HLT/NAACL*. Edmonton, Canada, pp. 127–133.
- Kurohashi, S., T. Nakazawa, K. Alexis, and D. Kawahara: 2005, 'Example-based Machine Translation Pursuing Fully Structural NLP'. In: *Proceedings of the International Workshop on Spoken Language Translation*. Pittsburgh, Pennsylvania, USA, pp. 207–212.
- Lepage, Y. and E. Denoual: 2006, 'The Purest EBMT System Ever Built'. *Machine Translation this volume*.
- Lin, D.: 2004, 'A Path-based Transfer Model for Machine Translation'. In: *Proceedings of International Conference on Computational Linguistics*. Geneva, Switzerland, pp. 625–630.
- Melamed, I. D.: 2004, 'Statistical Machine Translation by Parsing'. In: *Proceedings of the Annual Meeting of the Association for Computational Linguistics*. Barcelona, Spain, pp. 653–660.
- Menezes, A. and S. D. Richardson: 2003, 'A best-first alignment algorithm for extraction of transfer mappings'. In: M. Carl and A. Way (eds.): *Recent Advances in Example-Based Machine Translation*. Dordrecht, The Netherlands: Kluwer Academic Publishers, pp. 421–442.
- Och, F. and H. Ney: 2003, 'A Systematic Comparison of Various Statistical Alignment Models'. *Computational Linguistics* **29**(1), 19–51.
- Och, F. J.: 2003, 'Minimum Error Rate Training in Statistical Machine Translation'. In: *Proceedings of the Annual Meeting of the Association for Computational Linguistics*. Sapporo, Japan, pp. 160–167.
- Och, F. J., D. Gildea, S. Khudanpur, A. Sarkar, K. Yamada, A. Fraser, S. Kumar, L. Shen, D. Smith, K. Eng, V. Jain, Z. Jin, and D. Radev: 2004, 'A Smorgasbord of Features for

- Statistical Machine Translation'. In: *Proceedings of HLT/NAACL*. Boston, Massachusetts, USA, pp. 161–168.
- Och, F. J. and H. Ney: 2002, 'Discriminative Training and Maximum Entropy Models for Statistical Machine Translation'. In: *Proceedings of the Annual Meeting of the Association for Computational Linguistics*. Philadelphia, Pennsylvania, USA, pp. 295–302.
- Och, F. J. and H. Ney: 2004, 'The Alignment Template Approach to Statistical Machine Translation'. *Computational Linguistics* **30**(4), 417–449.
- Papineni, K., S. Roukos, T. Ward, and W.-J. Zhu: 2002, 'Bleu: a Method for Automatic Evaluation of Machine Translation'. In: *Proceedings of the Annual Meeting of the Association for Computational Linguistics*. Philadelphia, Pennsylvania, USA, pp. 311–318.
- Somers, H.: 2003, 'An Overview of EBMT'. In: M. Carl and A. Way (eds.): *Recent Advances in Example-Based Machine Translation*. Dordrecht, The Netherlands: Kluwer Academic Publishers, pp. 3–58.
- Vogel, S., Y. Zhang, F. Huang, A. Tribble, A. Venugopal, B. Zhao, and A. Waibel: 2003, 'The CMU Statistical Machine Translation System'. In: *MT Summit*. New Orleans, Louisiana, USA, pp. 402–409.
- Way, A. and N. Gough: 2005, 'Comparing Example-Based and Statistical Machine Translation'. *Natural Language Engineering* **11**(3), 295–309.
- Wu, D.: 1997, 'Stochastic Inversion Transduction Grammars and Bilingual Parsing of Parallel Corpora'. *Computational Linguistics* **23**(3), 377–403.
- Wu, D.: 2006, 'MT Model Space: Statistical vs. Compositional vs. Example-Based Machine Translation'. *Machine Translation* **this volume**.
- Yamada, K. and K. Knight: 2002, 'A Decoder for Syntax-Based Statistical MT'. In: *Proceedings of the Annual Meeting of the Association for Computational Linguistics*. Philadelphia, Pennsylvania, USA, pp. 303–310.

