

# Statistical Machine Translation Using Coercive Two-Level Syntactic Transduction

Charles Schafer and David Yarowsky

Center for Language and Speech Processing / Department of Computer Science  
Johns Hopkins University  
Baltimore, MD 21218 USA  
{cschafer,yarowsky}@cs.jhu.edu

## Abstract

We define, implement and evaluate a novel model for statistical machine translation, which is based on shallow syntactic analysis (part-of-speech tagging and phrase chunking) in both the source and target languages. It is able to model long-distance constituent motion and other syntactic phenomena without requiring a full parse in either language. We also examine aspects of lexical transfer, suggesting and exploring a concept of translation *coercion* across parts of speech, as well as a transfer model based on lemma-to-lemma translation probabilities, which holds promise for improving machine translation of low-density languages. Experiments are performed in both Arabic-to-English and French-to-English translation demonstrating the efficacy of the proposed techniques. Performance is automatically evaluated via the Bleu score metric.

## 1 Introduction

In this work we define, implement and evaluate a novel model for statistical machine translation (SMT).

Our goal was to produce a SMT system for translating foreign languages into English which utilizes some syntactic information in both the foreign language and English without, however, requiring a full parse in either language. Some advantages of not relying on full parses include that (1) there is a lack of availability of parsers for many languages of interest; (2) parsing time complexity represents a potential bottleneck for both model training and testing.

Intuitively, the explicit modeling of syntactic phenomena should be of benefit in the machine translation task; the ability to handle long-distance motion in an intelligently constrained way is a salient example of such a benefit. Allowing unconstrained translation reorderings at the word level generates a very large set of permutations that pose a difficult search problem at decoding time. We propose a model that makes use of shallow parses (text chunking) to support long-distance motion of phrases without requiring deeper analysis of syntax. The resources required to train this system on a new language are minimal, and we gain the ability to model long-distance movement and some interesting properties of lexical translation across parts of speech. One of

the source languages we examine in this paper, Arabic, has a canonical sentence-level order of Verb-Subject-Object, which means that translation into English (with a standard ordering of Subject-Verb-Object) commonly requires motion of entire phrasal constituents, which is not true of French-to-English translation, to cite one language pair whose characteristics have wielded great influence in the history of work on statistical machine translation. A key motivation for and objective of this work was to build a translation model and feature space to handle the above-described phenomenon effectively.

## 2 Prior Work

Statistical machine translation, as pioneered by IBM (e.g. Brown et al., 1993), is grounded in the noisy channel model. And similar to the related channel problems of speech and handwriting recognition, the original SMT language pair French-English exhibits a relatively close linear correlation in source and target sequence. Much common local motion that is observed for French, such as adjective-noun swapping, is adequately modeled by the relative-position-based distortion models of the classic IBM approach. Unfortunately, these distortion models are less effective for languages such as Japanese or Arabic, which have substantially different top-level sentential word orders from English, and hence longer distance constituent motion.

Wu (1997) and Jones and Havrilla (1998) have sought to more closely tie the allowed motion of constituents between languages to those syntactic transductions supported by the independent rotation of parse tree constituents. Yamada and Knight (2000, 2001) and Alshawi et al. (2000) have effectively extended such syntactic transduction models to fully functional SMT systems, based on channel model tree transducers and finite state head transducers respectively. While these models are well suited for the effective handling of highly divergent sentential word orders, the above frameworks have a limitation shared with probabilistic context free grammars that the preferred ordering of subtrees is insufficiently constrained by their embedding context, which is especially problematic for very deep syntactic parses.

In contrast, Och et al. (1999) have avoided the constraints of tree-based syntactic models and allow the rel-

atively flat motion of empirically derived phrasal chunks, which need not adhere to traditional constituent boundaries.

Our current paper takes a middle path, by grounding motion in syntactic transduction, but in a much flatter 2-level model of syntactic analysis, based on flat embedded noun-phrases in a flat sentential constituent-based chunk sequence that can be driven by syntactic brackets and POS tag models rather than a full parser, facilitating its transfer to lower density languages. The flatter 2-level structures also better support transductions conditioned to full sentential context than do deeply embedded tree models, while retaining the empirically observed advantages of translation ordering independence of noun-phrases.

Another improvement over Och et al. and Yamada and Knight is the use of the finite state machine (FSM) modelling framework (e.g. Bangalore and Riccardi, 2000), which offers the considerable advantage of a flexible framework for decoding, as well as a representation which is suitable for the fixed two-level phrasal modelling employed here.

Finally, the original cross-part-of-speech lexical coercion models presented in Section 4.3.3 have related work in the primarily-syntactic coercion models utilized by Dorr and Habash (2002) and Habash and Door (2003), although their induction and modelling are quite different from the approach here.

### 3 Resources

As in other SMT approaches, the primary training resource is a sentence-aligned parallel bilingual corpus. We further require that each side of the corpus be part-of-speech (POS) tagged and phrase chunked; our lab has previously developed techniques for rapid training of such tools (Cucerzan and Yarowsky, 2002). Our translation experiments were carried out on two languages: Arabic and French. The Arabic training corpus was a subset of the United Nations (UN) parallel corpus which is being made available by the Linguistic Data Consortium. For French-English training, we used a portion of the Canadian Hansards. Both corpora utilized sentence-level alignments publicly distributed by the Linguistic Data Consortium.

POS tagging and phrase chunking in English were done using the trained systems provided with the fnTBL Toolkit (Ngai and Florian, 2001); both were trained from the annotated Penn Treebank corpus (Marcus et al., 1993). French POS tagging was done using the trained French lexical tagger also provided with the fnTBL software. For Arabic, we used a colleague’s POS tagger and tokenizer (clitic separation was also performed prior to POS tagging), which was rapidly developed in our laboratory. Simple regular-expression-based phrase chunkers were developed by the authors for both Arabic and French, requiring less than a person-day each using existing multilingual learning tools.

A further input to our system is a set of word alignment links on the parallel corpus. These are used to compute word translation probabilities and phrasal alignments. The word alignments can in principle come from any source: a dictionary, a specialized alignment program, or another SMT system. We used alignments generated by Giza++ (Och and Ney, 2000) by running it in both directions (e.g., Arabic  $\rightarrow$  English **and** English  $\rightarrow$  Arabic) on our parallel corpora. The union of these bidirectional alignments was used to compute cross-language phrase correspondences by simple majority voting, and for purposes of estimating word translation probabilities, each link in this union was treated as an independent instance of word translation.

## 4 Translation Model

Now we turn to a detailed description of the proposed translation model. The exposition will give a formal specification and also will follow a running example throughout, using one of the actual Arabic test set sentences. This example, its gloss, system translation and reference human translation are shown in Table 1.

The translation model (TM) we describe is trained directly from counts in the data, and is a direct model, not a noisy channel model. It consists of three nested components: (1) a sentence-level model of phrase correspondence and reordering, (2) a model of intra-phrase translation, and (3) models of lexical transfer, or word translation. We make a key assumption in our construction that translation at each of these three levels is independent of the others.

### 4.1 Sentence Translation

As mentioned, both the foreign language and English corpora are input with “hard” phrase bracketings and labeled with “hard” phrase types (e.g., NP, VP<sup>1</sup>, PPNP<sup>2</sup>, etc.) as given by the output of the phrase chunker. These are denoted in the top-level model presentation in Table 2(1). Given word alignment links, as described in Section 2, we compute phrasal alignments on training data. We constrain these to have cardinality  $(foreign)N \rightarrow 1(English)$ . Next, we collect counts over aligned phrase sequences and use the relative frequencies to estimate the probability distribution in Table 2(2). Particularly for smaller training corpora, unseen foreign-language phrase sequences are a problem, so we implemented a simple backoff method which assigns probability to translations of unseen foreign-language phrase sequences. Table 2(3) encapsulates the remainder of the translation model, which is described below.

As an example, Table 3 shows the most probable aligned English phrase sequence generations given an Arabic simple sentence having the canonical VSO ordering. Also, note that all probabilities in the following

<sup>1</sup>VP in our parlance is perhaps more properly called a verb chunk: it consists of a verb, its auxiliaries, and contiguous adverbs.

<sup>2</sup>PPNP consists of a NP with its prepositional head attached.

Arabic Example Sentence From Test Set	
(ARABIC)	<i>twSy Al- ljmp Al- sAdsp Al- jmEyp Al- EAmp b- AEmAd m\$rwE Al- mqrr Al- tAly :</i>
(PHR.-BRACKETED AR.)	<i>[twSy] [Al- ljmp Al- sAdsp] [Al- jmEyp Al- EAmp] [b- AEmAd m\$rwE Al- mqrr Al- tAly] [:]</i>
(AN ENG. GLOSS)	[recommends] [the committee the sixth] [the assembly the general] [to adoption draft the decision the following] [:]
(ENG. MT OUTPUT)	[the sixth committee] [recommends] [the general assembly] [in the adoption of the following draft resolution] [:]
(REFERENCE TRANS.)	the sixth committee recommends to the general assembly the adoption of the following draft decision :

Table 1: An Arabic translation from the test set. We revisit portions of this example throughout the text. All Arabic strings in this paper are rendered in the reversible Buckwalter transliteration. In addition, all words or symbols referring to Arabic and French in this paper are italicized.

figures and tables are from the actual Arabic and French trained systems.

Arabic Phrase Sequence	Aligned English Phrase Sequence	Prob.
<i>VP<sub>1</sub> NP<sub>2</sub> NP<sub>3</sub></i>	<i>NP<sub>2</sub> VP<sub>1</sub> NP<sub>3</sub></i>	0.23
<i>VP<sub>1</sub> NP<sub>2</sub> NP<sub>3</sub></i>	<i>VP<sub>1</sub> NP<sub>2</sub> PP<sub>3</sub></i>	0.10
<i>VP<sub>1</sub> NP<sub>2</sub> NP<sub>3</sub></i>	<i>NP<sub>3</sub> VP<sub>1,2</sub></i>	0.06

Table 3: Top learned sentence-level reorderings for Arabic, for canonical Arabic simple sentence structure VP (verb) NP (subject) NP (object). Subscripts in English phrase sequence are alignments to positions in the corresponding Arabic phrase sequence.

## 4.2 Phrase Translation

Given an Arabic test sentence, a distribution of aligned English phrase sequences is proposed by the sentence-level model described in the previous section and in Table 2. Each proposed English phrase in each of the phrase sequence possibilities, therefore, comes to the middle level of the translation model with access to the identity of the French phrases aligned to it. Phrase translation is implemented as shown in Table 4. The phrase translation model is structured with several levels of backoff: if no observations exist from training data for a particular level, the model backs off to the next-more-general level. In all cases, generation of an English phrase is conditioned on the foreign phrase as well as the type (NP, VP, etc.) of the English phrase.

Table 4 (1) describes the initial phrase translation model. It comes into play if the precise sequence of foreign words has been observed aligning to an English phrase of the appropriate type. In the example, we are trying to generate an NP given the Arabic word string “*Al- ljmp Al- sAdsp*” (literally: “the committee the sixth”). If this has been observed in data, then that relative frequency distribution serves as the translation probability distribution. Table 11 contains examples of some of these literal phrase translations from the French data.

The next stage of backoff from the above, literal level is a model that generates aligned English POS tag sequences given foreign POS tag sequences: details and an example can be found in Table 4(2). The sequence alignments determine the position in English phrase and the part-of-speech into which we translate the foreign

word. Again, translation is also conditioned on the English phrase type. Table 5 and Table 6 show the most probable aligned English sequence generations for two of the phrases in the example sentence.

If there were no counts for (foreign-POS-sequence, english-phrase-type) then we back off to counts collected over (foreign-coarse-POS-sequence, english-phrase-type), where a coarse POS is, for example, *N* instead of *NOUN-SG*. This is shown in Table 4(3).

In case further backoff is needed, as shown in Table 4(4), we begin stripping POS-tags off the “less significant” (non-head) end of the foreign POS-sequence until we are left with a phrase sequence that has been seen in training, and from this a corresponding English phrase distribution is observable. We define the “less significant” end of a phrase to be the end if it is head-initial, or the beginning if it is head-final, and at this point ignore issues such as nested structure in French and Arabic NP’s.

Aligned English POS-tag Sequence Translation Probabilities (conditioned on Arabic POS-tag sequence from NP in example)	
$P(DT_{\emptyset} JJ_4 NN_2)$	$ DET_1 NOUN-SG_2 DET_3 ADJ_4, NP) = 0.22$
$P(JJ_4 NN_1)$	$ DET_1 NOUN-SG_2 DET_3 ADJ_4, NP) = 0.20$
$P(DT_{\emptyset} NN_1)$	$ DET_1 NOUN-SG_2 DET_3 ADJ_4, NP) = 0.13$
$P(DT_{\emptyset} VBN_4 NNS_2)$	$ DET_1 NOUN-SG_2 DET_3 ADJ_4, NP) = 0.13$
$P(DT_1 NN_2)$	$ DET_1 NOUN-SG_2 DET_3 ADJ_4, NP) = 0.04$
$P(DT_3 JJ_4 NN_2)$	$ DET_1 NOUN-SG_2 DET_3 ADJ_4, NP) = 0.03$
$P(DT_1 VBN_4 NNS_2)$	$ DET_1 NOUN-SG_2 DET_3 ADJ_4, NP) = 0.03$
$P(DT_{\emptyset} NN_4 NN_2)$	$ DET_1 NOUN-SG_2 DET_3 ADJ_4, NP) = 0.02$
$P(JJ_4 NNS_2)$	$ DET_1 NOUN-SG_2 DET_3 ADJ_4, NP) = 0.02$
$P(DT_1 JJ_4 NN_2)$	$ DET_1 NOUN-SG_2 DET_3 ADJ_4, NP) = 0.02$
$P(NN_2)$	$ DET_1 NOUN-SG_2 DET_3 ADJ_4, NP) = 0.02$
$P(NN_4 NN_2)$	$ DET_1 NOUN-SG_2 DET_3 ADJ_4, NP) = 0.02$

Table 5: From the running Arabic example, top English NP generations given an Arabic phrase *DET NOUN-SG DET ADJ*. Note:  $\emptyset$  denotes a null alignment (generation from null). Generation from a null alignment is allowed for specified parts of speech, such as determiners and prepositions.

## 4.3 Lexical Transfer

### 4.3.1 The Basic Model

In the basic model of word generation, phrases may be translated directly as single atomic entities (as in Table 4(1)), or via phrasal decomposition to individual words translated independently, conditioned only on the source word and target POS. Word translation in the latter case

Top-level Definition of Translation Model		
Example Instantiation of Model Variables		Model Description
$P(\text{the sixth committee recommends the general assembly ..}   \text{twSy Al- ljmp Al- sAdsp Al- jmEyp Al- EAmp ..}) =$		$P(\text{english\_words}   \text{foreign\_words}) =$
$P([\text{twSy}]_{VP_1} [\text{Al- ljmp Al- sAdsp}]_{NP_1} [\text{Al- jmEyp Al- EAmp}]_{NP_2} ..   \text{twSy Al- ljmp Al- sAdsp Al- jmEyp Al- EAmp ..})$	(1)	$P(\text{foreign\_bracketing} , \text{foreign\_phrase\_sequence}   \text{foreign\_words})$
$*P(NP_2 VP_1 NP_3 PPNP_4 PUNC_5   VP_1 NP_2 NP_3 PPNP_4 PUNC_5)$	(2)	$P(\text{english\_phrase\_sequence} , \text{phrase\_alignment\_matrix}   \text{foreign\_phrase\_sequence})$
$*P([\text{the sixth committee}]_{NP_2} [\text{recommends}]_{VP_1} [\text{the general assembly}]_{NP_3} ..   [\text{twSy}]_{VP_1} [\text{Al- ljmp Al- sAdsp}]_{NP_1} [\text{Al- jmEyp Al- EAmp}]_{NP_2} .. , NP_2 VP_1 NP_3 PPNP_4 PUNC_5)$	(3)	$P(\text{english\_words} , \text{english\_bracketing} , \text{english\_phrase\_sequence}   \text{foreign\_words} , \text{foreign\_bracketing} , \text{foreign\_phrase\_sequence} , \text{english\_phrase\_sequence} , \text{phrase\_alignment\_matrix})$

Table 2: Statement of the translation model at top level.

Phrase Translation Model with Backoff Pathways		
Example Instantiations		Model Statement
$P(\text{the sixth committee}   \text{Al- ljmp Al- sAdsp} , \text{NP}) =$		
$P(\text{the sixth committee}   \text{Al- ljmp Al- sAdsp} , \text{NP})$	(1)	$P(W_{E_1} W_{E_2} .. W_{E_n}   W_{F_1} W_{F_2} .. W_{F_m} , \text{phr\_type}_E)$
$\downarrow$		$\downarrow (\text{backoff if } C(W_{F_1} W_{F_2} .. W_{F_m} , \text{phr\_type}_E) = 0)$
$P(DT_1 JJ_4 NN_2   DET_1 NOUN\text{-}SG_2 DET_3 ADJ_4 , \text{NP})$	(2)	$P(T_{fine_{E_1}} T_{fine_{E_2}} .. T_{fine_{E_n}} , \Xi_i   T_{fine_{F_1}} T_{fine_{F_2}} .. T_{fine_{F_m}} , \text{phr\_type}_E)$
$*P(\text{the}   \text{Al-} , \text{DT})$		$*P(W_{E_1}   W_{F_{\Xi_i(1)}} , T_{fine_{E_1}})$
$*P(\text{committee}   \text{ljmp} , \text{NN})$		$*P(W_{E_2}   W_{F_{\Xi_i(2)}} , T_{fine_{E_2}})$
$*P(\text{sixth}   \text{sAdsp} , \text{JJ})$		$*.. * P(W_{E_n}   W_{F_{\Xi_i(n)}} , T_{fine_{E_n}})$
$\downarrow$		$\downarrow (\text{backoff if } C(T_{fine_{F_1}} T_{fine_{F_2}} .. T_{fine_{F_m}} , \text{phr\_type}_E) = 0)$
$P(DT_1 JJ_4 NN_2   D_1 N_2 D_3 A_4 , \text{NP})$	(3)	$P(T_{fine_{E_1}} T_{fine_{E_2}} .. T_{fine_{E_n}} , \Xi_i   T_{coarse_{F_1}} T_{coarse_{F_2}} .. T_{coarse_{F_m}} , \text{phr\_type}_E)$
$*P(\text{the}   \text{Al-} , \text{DT})$		$*P(W_{E_1}   W_{F_{\Xi_i(1)}} , T_{fine_{E_1}})$
$*P(\text{committee}   \text{ljmp} , \text{NN})$		$*P(W_{E_2}   W_{F_{\Xi_i(2)}} , T_{fine_{E_2}})$
$*P(\text{sixth}   \text{sAdsp} , \text{JJ})$		$*.. * P(W_{E_n}   W_{F_{\Xi_i(n)}} , T_{fine_{E_n}})$
$\downarrow$		$\downarrow (\text{backoff if } C(T_{coarse_{F_1}} T_{coarse_{F_2}} .. T_{coarse_{F_m}} , \text{phr\_type}_E) = 0)$
$P(?   D_1 N_2 D_3 , \text{NP})$	(4)	$P(T_{fine_{E_1}} T_{fine_{E_2}} .. T_{fine_{E_n}} , \Xi_i   T_{coarse_{F_1}} T_{coarse_{F_2}} .. T_{coarse_{F_{m-1}}} , \text{phr\_type}_E)$
$*? * .. * ?$		$*? * .. * ?$
$\downarrow$		$\downarrow (\text{backoff if } C(T_{coarse_{F_1}} T_{coarse_{F_2}} .. T_{coarse_{F_{m-1}}} , \text{phr\_type}_E) = 0)$
$P(?   D_1 N_2 , \text{NP})$	(4)	$P(T_{fine_{E_1}} T_{fine_{E_2}} .. T_{fine_{E_n}} , \Xi_i   T_{coarse_{F_1}} T_{coarse_{F_2}} .. T_{coarse_{F_{m-2}}} , \text{phr\_type}_E)$
$*? * .. * ?$		$*? * .. * ?$
$\downarrow$		$\downarrow (\text{backoff if } C(T_{coarse_{F_1}} T_{coarse_{F_2}} .. T_{coarse_{F_{m-2}}} , \text{phr\_type}_E) = 0)$
....		....

Table 4: The phrase translation model, with backoff. Examples on the left side are from one of the Arabic test sentences. (1) is the direct, lexical translation level. (2) - (4) constitute the backoff path to handle detailed phenomena unseen in the training set. (2) is a model of fine POS-tag reordering and lexical generation; (3) is similar, but conditions generation on *coarse* POS-tag sequences in the foreign language. (4) is a model for progressively stripping off POS-tags from the “less significant” end of a foreign sequence. The idea is to do this until we reach a subsequence that has been seen in training data, and which we therefore have a distribution of valid generators for. The term  $\Xi_i$  in (2) - (4) is a position alignment matrix. At all times, we generate not just an English POS-tag sequence, but rather an **aligned** sequence. Similarly, in the lexical transfer probabilities shown in this table, there is a function  $\Xi_i()$  which takes an English sequence position index and returns the (unique) foreign word position to which it is aligned<sup>4</sup>.

Aligned English POS-tag Sequence Translation Probabilities (conditioned on Arabic POS-tag sequence from VP in example)	
$P(\text{VBZ}_1   \text{VERB-IMP}_1 , \text{VP}) = .28$	
$P(\text{VBP}_1   \text{VERB-IMP}_1 , \text{VP}) = .17$	
$P(\text{VBD}_1   \text{VERB-IMP}_1 , \text{VP}) = .09$	
....	
$P(\text{MD}_0 \text{VB}_1   \text{VERB-IMP}_1 , \text{VP}) = .06$	

Table 6: From the Arabic example, top English VP generations given an Arabic phrase *VERB-IMP*.

is done in the context that the model has already proposed a sequence of POS tags for the phrase. Thus we

know the English POS of the word we are trying to generate in addition to the foreign word that is generating it. Consequently, we condition translation on English POS as well as the foreign word. Table 7 describes the backoff path for basic lexical transfer and presents a motivating example in the French word *droit*. Translation probabilities for one of the words in the example Arabic sentence can be found in Table 8.

#### 4.3.2 Generation via a Lemma Model

To counter sparse data problems in estimating word translation probabilities, we also implemented a lemma-

Word Generation		Model with Backoff Pathways
Examples		
$P(W_E   \textit{droit}, \text{NNS})$		$P(W_E   W_F, T_{fine_E})$
rights	0.4389	$p(\textit{rights}   \textit{droit}, \text{NNS})$
benefits	0.0690	
people	0.0533	
laws	0.0188	
		$\downarrow$ (backoff if $C(W_F, T_{fine_E}) = 0$ )
$P(W_E   \textit{droit}, \text{N})$		$P(W_E   W_F, T_{coarse_E})$
right	0.4970	
law	0.1318	
rights	0.0424	$p(\textit{rights}   \textit{droit}, \text{N})$
property	0.0115	
		$\downarrow$ (backoff if $C(W_F, T_{coarse_E}) = 0$ )
$P(W_E   \textit{droit})$		$P(W_E   W_F)$
right	0.2919	
entitled	0.0663	
law	0.0652	
the	0.0249	
to	0.0240	
rights	0.0210	$p(\textit{rights}   \textit{droit})$
		$\downarrow$ (backoff if $C(W_F) = 0$ )
		$p(\text{UNKNOWN\_WORD}   W_F) = 1$

Table 7: Description of the conditioning for different levels of backoff in the lexical transfer model. The example shows translations for the French word *droit* (“right”) conditioned on decreasingly specific values. The progressively lower ranking of the correct translation as we move from fine, to coarse, to no POS, illustrates the benefit of conditioning generation on the English part of speech.

Arabic Word	English POS	English Wd.	Prob.
<i>ljnp</i>	NN	committee	0.591
<i>ljnp</i>	NN	commission	0.233
<i>ljnp</i>	NN	subcommittee	0.035
<i>ljnp</i>	NN	acc	0.013
<i>ljnp</i>	NN	report	0.005
<i>ljnp</i>	NN	ece	0.004
<i>ljnp</i>	NN	icrc	0.004
<i>ljnp</i>	NN	aalcc	0.004
<i>ljnp</i>	NN	escap	0.004
<i>ljnp</i>	NN	escwa	0.004
<i>ljnp</i>	NN	eca	0.003
<i>ljnp</i>	NNS	members	0.088
<i>ljnp</i>	NNS	recommendations	0.033
<i>ljnp</i>	NNS	copuos	0.033
<i>ljnp</i>	NNS	questions	0.027
<i>ljnp</i>	NNS	representatives	0.024
<i>ljnp</i>	N	committee	0.577
<i>ljnp</i>	N	commission	0.227
<i>ljnp</i>	N	subcommittee	0.035

Table 8: From running example, translation probabilities for Arabic noun *ljnp*, “committee”.

based model for word translation. Under this model, translation distributions are estimated by counting word alignment links between foreign and English lemmas, assuming a lemmatization of both sides of the parallel corpus as input. The form of the model is illustrated below:

$$\begin{aligned}
 P(W_E | W_F, T_{coarse_F}, T_{fine_E}) = & \\
 & P(W_E | \textit{lemma}_E, T_{coarse_F}, T_{fine_E}) * \\
 & P(\textit{lemma}_E | \textit{lemma}_F, T_{coarse_F}, T_{fine_E}) * \\
 & P(\textit{lemma}_F | W_F, T_{coarse_F}, T_{fine_E}) \\
 & \downarrow \text{approximated by} \\
 & P(W_E | \textit{lemma}_E, T_{fine_E}) * \\
 & P(\textit{lemma}_E | \textit{lemma}_F, T_{coarse_E}) * \\
 & P(\textit{lemma}_F | W_F, T_{coarse_F})
 \end{aligned}$$

First, note that  $P(\textit{lemma}_F | W_F, T_{coarse_F})$  is very simply a hard lemma assignment by the foreign language lemmatizer. Second, English word generation from English lemma and coarse POS ( $P(W_E | \textit{lemma}_E, T_{fine_E})$ ) is programmatic, and can be handled by means of rules in conjunction with a lookup table for irregular forms. The only distribution here that must be estimated from data is  $P(\textit{lemma}_E | \textit{lemma}_F, T_{coarse_E})$ . This is done as described above. Furthermore, given an electronic translation dictionary, even this distribution can be pre-loaded: indeed, we expect this to be an advantage of the lemma model, and an example of a good opportunity for integrating compiled human knowledge about language into an SMT system. Some examples of the lemma model combating sparse data problems inherent in the basic word-to-word models can be found in Table 9.

### 4.3.3 Coercion

Lexical coercion is a phenomenon that sometimes occurs when we condition translation of a foreign word on the word and the English part-of-speech. We find that the system we have described frequently learns this behavior: specifically, the model learns in some cases how to generate, for instance, a nominal form with similar meaning from a French adjective, or an adjectival realization of a French verb’s meaning; some examples of this phenomenon are shown in Table 10. We find this coercion effect to be of interest because it identifies interesting associations of meaning. For example, in Table 10 “willing” and “ready” are both sensible ways to realize the meaning of the action “to accept” in a passive, descriptive mode. *droit* behaves similarly. Though the English verb “to right” or “to be righted” does not have the philosophical/judicial entitlement sense of the noun “right”, we see that the model has learned to realize the meaning in an active, verbal form: e.g., VBG ‘receiving’ and VB “qualify”.

## 5 Decoding

Decoding was implemented by constructing finite-state machines (FSMs) **per evaluation sentence** to encode relevant portions (for the individual sentence in question) of the component translation distributions described above. Operations on these FSMs are performed using the AT&T FSM Toolkit (Mohri et al., 1997). The FSM constructed for a test sentence is subsequently composed with a FSM trigram language model created via the SRI Language Modeling Toolkit (Stolcke, 2002). Thus we use the trigram language model to implement rescoring

of the (direct) translation probabilities for the English word sequences in the translation model lattice.

We found that using the finite-state framework and the general-purpose AT&T toolkit greatly facilitates decoder development by freeing the implementation from details of machine composition and best-path searching, etc.

The structure of the translation model finite-state machines is as illustrated in Figure 1. The sentence-level (aligned phrase sequence generation) and phrase-level (aligned intra-phrase sequence generation) translation probabilities are encoded on epsilon arcs in the machines. Word translation probabilities are placed onto arcs emitting the word as an output symbol (in the figure, note the arcs emitting “committee”, “the”, etc.). The FSM in Figure 1 corresponds to the Arabic example sentence used throughout this paper. In the portion of the machine shown, the (best) path which generated the example sentence is drawn in bold. Finally, Figure 2 is a rendering of the actual FSM (aggressively pruned for display purposes) that generated the example Arabic sentence; although labels and details are not visible, it may provide a visual aid for better understanding the structure of the FSM lattices generated here.

As a practical matter in decoding, during translation model FSM construction we modified arc costs for output words in the following way: a fixed bonus was assigned for generating a “content” word translating to a “content” word. Determining what qualifies as a content word was done on the basis of a list of content POS tags for each language. For example, all types of nouns, verbs and adjectives were listed as content tags; determiners, prepositions, and most other closed-class parts of speech were not. This implements a reasonable penalty on undesirable output sentence lengths. Without such a penalty, translation outputs tend to be very short: long sentence hypotheses are penalized *de facto* merely by containing many word translation probabilities. An additional trick in decoding is to use only the N-best translation options for sentence-level, phrase-level, and word-level translation. We found empirically (and very consistently) in dev-test experiments that restricting the syntactic transductions to a 30-best list and word translations to a 15-best list had no negative impact on Bleu score. The benefit, of course, is that the translation lattices are dramatically reduced in size, speeding up composition and search operations.

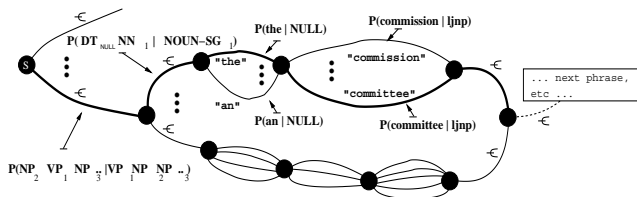


Figure 1: An illustration of the translation model structure for an Arabic test sentence.

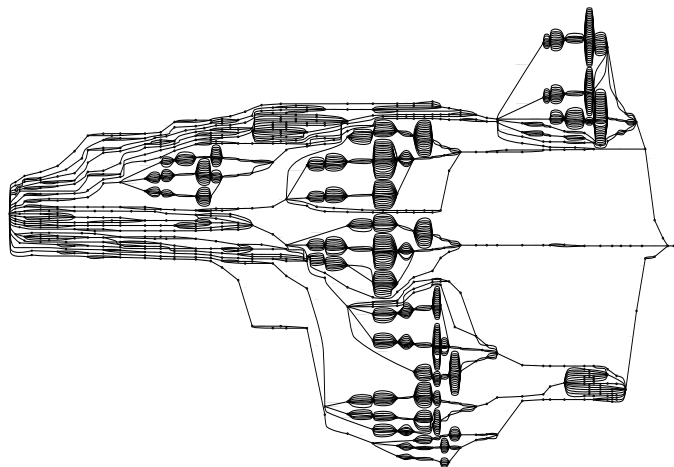


Figure 2: A portion of the translation model for an Arabic test sentence, compacted and aggressively pruned by path probability for display purposes.

## 6 Evaluation

Results Tables A and B below list evaluation results for translation on the Arabic and French test sets respectively. In each case, results for a comparison system – the Giza++ IBM Model 4 implementation (Och and Ney, 2000) with the ReWrite decoder (Marcu and Germann, 2002) – are included as a benchmark. Results were generated for training corpora of varying sizes. For Arabic, we ran our system on two large subsets of the UN corpus and evaluated on a 200-sentence held-out set (refer to Results Table A below). For the 150K sentence Arabic training set, Giza++ and the shallow syntax model achieved very similar performance. We were unable to obtain evaluation numbers for Giza++/ReWrite on the large Arabic training set, however, since its language model component has a vocabulary size limit which was exceeded in the larger corpus. In French we observed the systems to perform similarly on the small training sets we used (Results Table B). We performed some experiments in classifier combination using the two compatible (150K-training-sentence) Arabic systems, wherein a small devtest set was used to identify simple system combination parameters based on model confidence and sentence length. In situations where our system was confident we used its output, and used Giza++ output otherwise. We achieved a 3% boost in Bleu score over Giza++ performance on the evaluation set with these very simple classifier combination techniques, and anticipate that research in this direction – classifier combination of diversely trained SMT systems – could yield significant performance improvements.

System	Bleu Score	
	150K Trn. Sent.	500K Trn. Sent.
Giza++/ReWrite Decoder	0.17	*
2-level Syntax Model	0.17	0.18

Results Table A: Results comparison for Arabic to English translation on the UN corpus, with a 200-sentence evaluation set. Note that Giza++/ReWrite cannot be run for the 500K sentence training set; the CMU Language Modeling Toolkit, which ReWrite uses, has a vocabulary size limit which is exceeded in the 500K corpus.

System	Bleu Score	
	5K Trn. Sent.	20K Trn. Sent.
Giza++/ReWrite Decoder	0.08	0.11
2-level Syntax Model	0.08	0.09

Results Table B: Results comparison for French to English translation on the Canadian Hansards corpus (200-sentence evaluation set).

## 7 Conclusions

We have described and implemented an original syntax-based statistical translation model that yields baseline results which compete successfully with other state-of-the-art SMT models. This is particularly encouraging in that the authors are not well-versed in Arabic or French and it appears that the quality of the rule-based phrase chunkers we developed in a single person-day offers substantial room for improvement. We expect to be able to attain improved bracketings from native speakers and, in addition, via translational projection of existing bracketers. Secondly, the lemma model we have proposed for lexical transfer provides an efficient framework for integrating electronic dictionaries into SMT models. Although we have at this time no large electronic dictionaries for either Arabic or French, efforts are underway to acquire electronic or scanned paper dictionaries for this purpose. We did evaluate the lemma models in isolation for French and Arabic without dictionary inclusion, but in each experiment the results did not differ significantly from the word-specific lexical transfer models, despite their substantially reduced dimensionality. We anticipate that the relatively seamless direct incorporation of dictionaries into the lemma-based models will be particularly effective for translating low-density languages, which suffer from data sparseness in the face of limited parallel text. Finally, we incorporated lexical translation coercion models into this full SMT framework, the induction of which is a phenomenon of interest in its own right.

## 8 Acknowledgements

This work was supported in part by NSF grant number IIS-9985033. In addition, we owe many thanks to colleagues who generously lent their time and insights. David Smith shared his tools for Arabic part-of-speech tagging and morphological analysis and answered many questions about the Arabic language. Thanks to Skankar Kumar and Sanjeev Khudanpur for numerous helpful discussions.

## 9 References

- H. Alshawi, S. Bangalore, and S. Douglas. 2000. Learning dependency translation models as collections of finite state head transducers *Computational Linguistics*, 26(1), 45–60.
- S. Bangalore and G. Riccardi. 2000. Stochastic finite-state models for spoken language machine translation. In *Proceedings of the Workshop on Embedded Machine Translation Systems.*, pp. 52–59.
- P. Brown, S. Della Pietra, V. Della Pietra and R. Mercer. 1993. The mathematics of statistical machine translation: Parameter estimation. *Computational Linguistics*, 12(2), 263–312.
- S. Cucerzan and D. Yarowsky. 2002. Bootstrapping a Multilingual Part-of-speech Tagger in One Person-day. *Proceedings of the Sixth Conference on Natural Language Learning (CoNLL)*, Taipei, 2002.
- B. Dorr and N. Habash. 2002. Interlingua approximation: A generation-heavy approach. In *Proceedings of AMTA-2002*.
- W. A. Gale and K. W. Church. 1991. A Program for Aligning Sentences in Bilingual Corpora. In *29th Annual Meeting of the ACL*, Berkeley, CA.
- N. Habash and B. Dorr. 2003. A categorial variation database for English. In *Proceedings of NAACL-HLT 2003*
- D. Jones and R. Havrilla. 1998. Twisted pair grammar: Support for rapid development of machine translation for low density languages. In *Proceedings of AMTA98*, pp. 318–332.
- D. Marcu and U. Germann. 2002. *The ISI ReWrite Decoder Release 0.7.0b*. <http://www.isi.edu/licensed-sw/rewrite-decoder/>.
- M. Marcus, B. Santorini, and M. Marcinkiewicz. 1993. Building a large annotated corpus of English: the Penn Treebank. *Computational Linguistics*, Vol. 19.
- M. Mohri, F. Pereira, and M. Riley. 1997. ATT General-purpose finite-state machine software tools. <http://www.research.att.com/sw/tools/fsm/>.
- G. Ngai and R. Florian. Transformation-based learning in the fast lane. In *Proceedings of North American ACL 2001*, pages 40–47, June 2001.
- F. J. Och and H. Ney. 2000. Improved statistical alignment models. In *Proceedings of the 38th Annual Meeting of the Association for Computational Linguistics*, pages 440–447.
- F.J. Och, C. Tillmann, H. Ney. Improved Alignment Models for Statistical Machine Translation. In *Proceedings of EMNLP 1999*, pp. 20–28.
- K. Papineni, S. Roukos, T. Ward, and W. Zhu. 2001. Bleu: a method for automatic evaluation of machine translation. Technical Report RC22176 (W0109-022), IBM Research Division.
- A. Stolcke. 2002. SRILM - an extensible language modeling toolkit. In *Proceedings of the International Conference on Spoken Language Processing*, pages 901–904. Denver, CO, USA. <http://www.speech.sri.com/projects/srilm/>.
- D. Wu. 1997. Stochastic inversion transduction grammars and bilingual parsing of parallel corpora. *Computational Linguistics*, 23(3), 377–404.
- K. Yamada and K. Knight. 2001. A syntax-based statistical translation model. In *Proceedings of ACL-2001*, pp. 523–529.
- K. Yamada and K. Night. 2002. A decoder for syntax-based statistical MT In *Proceedings of ACL-2002*, pp. 303–310.

Word Translation Probabilities			
Word translation for <i>mangeait</i> conditioned on French Word, English POS			
<i>mangeait</i>	VBG	eating	1.00
<i>mangeait</i>	VB	go	0.50
<i>mangeait</i>	VB	anticipate	0.50
<i>mangeait</i>	VBD	were	1.00
<i>mangeait</i>	VBP	knelt	1.00
<i>mangeait</i>	NN	bill	1.00
Word translation for <i>mangeait</i> conditioned on French Word, English Coarse POS			
<i>mangeait</i>	V	eating	0.44
<i>mangeait</i>	V	were	0.22
<i>mangeait</i>	V	knelt	0.11
<i>mangeait</i>	V	go	0.11
<i>mangeait</i>	V	anticipate	0.11
<i>mangeait</i>	N	bill	1.00
Word translation for <i>mangeait</i> conditioned on French Word only			
<i>mangeait</i>		eating	0.29
<i>mangeait</i>		were	0.14
<i>mangeait</i>		go	0.07
<i>mangeait</i>		bill	0.07
Word translation for <i>mangeant</i> conditioned on French Word, English POS			
<i>mangeant</i>	RB	mostly	1.00
<i>mangeant</i>	JJ	fi nal	1.00
<i>mangeant</i>	VBN	obtained	1.00
<i>mangeant</i>	VBG	eating	1.00
<i>mangeant</i>	WP	who	1.00
<i>mangeant</i>	IN	through	1.00
<i>mangeant</i>	NN	lard	1.00
<i>mangeant</i>	VBZ	eats	0.50
<i>mangeant</i>	VBZ	comes	0.50
Lemma Translation Probabilities			
Generation of a verb lemma given <i>manger</i>			
<i>manger</i>	V	eat	0.60
<i>manger</i>	V	feed	0.05
<i>manger</i>	V	have	0.04
Generation of a noun lemma given <i>manger</i>			
<i>manger</i>	N	meal	0.06
<i>manger</i>	N	trough	0.04
<i>manger</i>	N	loaf	0.04
<i>manger</i>	N	food	0.04
Generation of an adj. lemma given <i>manger</i>			
<i>manger</i>	J	hungry	0.33
Raw lemma translation probabilities (ignoring English Coarse POS)			
<i>manger</i>		eat	0.28
<i>manger</i>		to	0.03
<i>manger</i>		feed	0.03
<i>manger</i>		out	0.02
<i>manger</i>		have	0.02
<i>manger</i>		are	0.02
<i>manger</i>		,	0.02
<i>manger</i>		you	0.01
<i>manger</i>		meal	0.01

Table 9: **Direct generation** (word-to-word translation probabilities at the various levels of backoff) is contrasted with **lemma generation**. *Manger* (“to eat”) is a relatively rare word in the Hansards. Note that due to low counts, the desired verb POS (target of generation) for ‘eat’ may not have been observed as a translation in training data. In addition, in this situation, noisy word alignments may cause an incorrect translation to have similar estimated translation probability. This problem is addressed by the lemma model; note the much sharper probability distribution for verb lemmas given *manger*. Generation of English inflections given lemma and target POS is algorithmic (and irregular exceptions are handled via a lookup table).

French Wd.	Eng. POS	Eng. Wd.	Prob.
<i>accepter</i>	JJ	unacceptable	0.12
<i>accepter</i>	JJ	acceptable	0.12
<i>accepter</i>	JJ	willing	0.11
<i>accepter</i>	JJ	ready	0.03
<i>accepter</i>	NN	acceptance	0.09
<i>accepter</i>	NN	amendment	0.03
<i>droit</i>	VBN	entitled	0.66
<i>droit</i>	VBN	allowed	0.09
<i>droit</i>	VBN	denied	0.03
<i>droit</i>	VBN	given	0.02
<i>droit</i>	VBN	permitted	0.02
<i>droit</i>	VBN	justifi ed	0.01
<i>droit</i>	VBN	qualifi ed	0.01
<i>droit</i>	VBN	allotted	0.01
<i>droit</i>	VB	qualify	0.14
<i>droit</i>	VB	be	0.11
<i>droit</i>	VB	have	0.09
<i>droit</i>	VB	receive	0.08
<i>droit</i>	VB	get	0.07
<i>droit</i>	VB	expect	0.03
<i>droit</i>	VBG	receiving	0.11
<i>droit</i>	VBG	getting	0.08
<i>droit</i>	NNS	rights	0.44
<i>droit</i>	NNS	benefi ts	0.69

Table 10: Examples of word translation coercions. Coercions of the French verb *accepter* “to accept” and the French noun *droit* “right” (there is parallel polysemy between the two languages for this word, but the predominant sense in our corpus is the philosophical/judicial sense, as opposed to the direction).

Eng. Phrase Type	French Phrase	Eng. Phrase	Prob.
NP	<i>dans_Le_cas_présent</i>	a_situation	0.25
NP	<i>dans_Le_cas_présent</i>	the_subject_of_debate	0.25
NP	<i>dans_Le_cas_présent</i>	the_position	0.25
NP	<i>dans_Le_cas_présent</i>	it	0.25
VP	<i>dans_Le_cas_présent</i>	should_apply	1.00
ADVP	<i>dans_Le_cas_présent</i>	really	1.00
PPNP	<i>dans_Le_cas_présent</i>	in_this_case	0.63
PPNP	<i>dans_Le_cas_présent</i>	in_this_instance	0.04
PPNP	<i>dans_Le_cas_présent</i>	in_this_actual_case	0.04
PPNP	<i>dans_Le_cas_présent</i>	in_this_particular_case	0.04
PPNP	<i>dans_Le_cas_présent</i>	in_that_case	0.04
PPNP	<i>dans_Le_cas_présent</i>	in_the_present_circumstances	0.04
VP	<i>acceptons</i>	accept	0.48
VP	<i>acceptons</i>	agree	0.14
NP	<i>acceptons</i>	this_consent	1.00
PPNP	<i>par_an</i>	per_year	0.67
PPNP	<i>par_an</i>	in_each_year	0.03
PPNP	<i>par_an</i>	for_a_year	0.03
ADVP	<i>par_an</i>	annually	1.00
NP	<i>par_an</i>	a_year	0.79
NP	<i>par_an</i>	each_year	0.02
NP	<i>un_discours</i>	a_speech	0.83
NP	<i>un_discours</i>	an_address	0.05
VP	<i>un_discours</i>	to_speak	1.00

Table 11: Examples of direct phrase translations (see Table 4(1)), including some coercions.