

Statistical Machine Translation: A Guide for Linguists and Translators

Mary Hearne and Andy Way*
 School of Computing, Dublin City University

Abstract

This paper presents an overview of Statistical Machine Translation (SMT), which is currently the dominant approach in Machine Translation (MT) research. In Way and Hearne (2010), we describe how central linguists and translators are to the MT process, so that SMT developers and researchers may better understand how to include these groups in continuing to advance the state-of-the-art. If these constituencies are to make an impact in the field of MT, they need to know how their input is used by SMT systems. Accordingly, our objective in this paper is to present the basic principles underpinning SMT in a way that linguists and translators will find accessible and useful.

1. Introduction

Corpus-based approaches to Machine Translation (MT) are founded on the use of parallel corpora. That is, previously unseen texts are automatically translated using information gleaned from examples of past translations produced by humans. The two primary MT strategies that belong to this paradigm are Example-Based MT (EBMT) and Statistical MT (SMT). While EBMT systems translate using the notion of ‘similarity’ to previously seen sentences, SMT systems do not. In this paper, we focus exclusively on a presentation of the basic principles underpinning the SMT approach. For a description of how the ideas behind SMT developed and how they relate to other MT approaches, and for a discussion of the central role played by translators and translations in SMT, see Way and Hearne (2010).

Statistical MT employs two distinct and separate processes: *training* and *decoding*. The training phase involves extracting a statistical model of translation from a parallel corpus, and a statistical model of the target language from a (typically much larger) monolingual corpus (Brown et al. 1990, 1993). The translation model effectively comprises a bilingual dictionary where each possible translation for a given source word or phrase has a probability associated with it. However, the model does not resemble a conventional dictionary where plausible entries only are permitted; many of the entries represent translations that are unlikely but not impossible, and the associated probabilities reflect this. The language model comprises a database of target-language word sequences (usually ranging between 1 and 7 words in length), each of which is also associated with a probability. Additional information beyond the language and translation models can be extracted during training, such as models of relative sentence length, word reordering, relative importance of translation vs. language model, etc. These induced models are then used during decoding, the process which actually yields a translation. The decoding process essentially treats translation as a search problem: given the sentence to be translated, search over all possible translations permitted by the translation model, and all possible reorderings thereof, for

CE: Gayathri		No. of pages: 21	
PE: Sharanya		Dispatch: 28.2.11	
		ToC head: [CM]	
		B	
		4	
		7	
		2	
		3	
Journal Name		Manuscript No.	
L N C		3	
			

1 the one which is assigned the highest overall probability according to the translation and
2 language models.

3 The idea of generating target sentences by translating words and phrases from the
4 source sentence in a random order using a model containing many non-sensical transla-
5 tions may not seem plausible. In fact, the methods used are not intended (in our opinion,
6 at least) to be either linguistically or cognitively plausible: see Way and Hearne (2010,
7 section 2.4) for further discussion of this topic. SMT is, however, *probabilistically* plausible;
8 rather than focusing on the best process to use to generate a single optimal translation for
9 a source sentence, SMT focuses on generating many thousands of hypothetical transla-
10 tions for the input string, and then working out which one of those is most likely.

11 While there are many variations of SMT systems, they all rely on these core principles
12 and techniques. Thus, SMT is not a generic term – the generic terms are corpus-based,
13 data-driven or perhaps probabilistic MT – but rather denotes a specific approach and
14 architecture. In this paper, we detail the main components that make up an SMT system,
15 covering the core training, parameter optimisation and decoding strategies. In Section 2,
16 we focus on the big picture, while in subsequent sections we zoom in on the finer
17 details: training in Section 3, including tuning in Section 3.3, and decoding in Section 4.
18 Section 5 provides a brief conclusion, together with some pointers to other accounts of
19 SMT for the interested reader.

20 2. Overview

21 At a high level, SMT gives us a view of MT expressed in a single formula. From this
22 vantage point, how translations are generated is irrelevant. The only salient issue is that
23 we can determine how likely any proposed translation is given the input string, and that
24 we can consequently determine the most probable (i.e. ‘best’, according to the system)
25 translation from a set of proposed candidates.

26 That is, SMT is founded on a well-defined decision problem: which translation is the
27 most likely? There are two formulae available to us to compute that score: the noisy-
28 channel model in (1), and the log-linear model in (2). The noisy-channel model is the
29 one traditionally used in the literature (Brown et al. 1990, 1993). While the log-linear
30 model can be instantiated to express precisely the same computation as the noisy-channel
31 one (as we explain later in this section), it is more flexible and has come into widespread
32 use in recent years (Och and Ney 2002).

$$33 \text{ Noisy-channel model: } Translation = \operatorname{argmax}_T P(S|T) \cdot P(T) \quad (1)$$

$$34 \text{ Log-linear model: } Translation = \operatorname{argmax}_T \sum_{m=1}^M \lambda_m \cdot h_m(T, S) \quad (2)$$

35 Firstly, we note that the leftmost parts of these equations ($Translation = \operatorname{argmax}_T \dots$) are
36 exactly the same and are to be interpreted as follows: assuming that we have T candidate
37 translations, we will let the output $Translation$ be the T with the maximum (argmax) score.
38 Secondly, we note that the rightmost part of the equations then differ in expressing how
39 each candidate translation should be scored.

40 The noisy-channel model comprises two component (or ‘feature’) scores $P(S|T)$ and
41 $P(T)$ which are to be multiplied together (as shown in (1)). The first of these, $P(S|T)$,
42 gives the likelihood that the source sentence S and the candidate translation T are transla-
43 tionally equivalent, i.e. that the meaning expressed in S is also captured in T . This feature
44
45
46
47
48
49

is generally referred to as the *translation model*. The second feature, $P(T)$, gives the likelihood that the candidate translation T is actually a valid sentence in the target language and is generally referred to as the *language model*.

Source Sentence:	Le chat entre dans la chambre.	
Adequate Fluent translation:	The cat enters the room.	
Adequate Disfluent translation:	The cat enters in the bedroom.	(3)
Fluent Inadequate translation:	My Granny plays the piano.	
Disfluent Inadequate translation:	piano Granny the plays My	

Thus, the traditional notions of adequacy and fluency used in human evaluations are modelled completely separately. As an example of how these might be used in practice, see the French-to-English examples in (3). Note that a translation would need to receive high adequacy and fluency scores in order to be deemed a good translation overall.

The translation model score is based on the extent to which the source sentence meaning is also expressed in the candidate translation; the details of how this score is computed are given in Section 3.2, but it is generally based on *lexical* correspondences only. Meanwhile, the language model score is based on the extent to which the candidate translation is likely to be a valid sentence in the target language regardless of whether or not its meaning bears any relationship to the source sentence. Details of how this score is computed are given in Section 3.1, but it is generally based on the occurrence frequencies of substrings of the sentence in target-language corpora. However, the final score depends on *both* of these features and the objective is to find the translation with the best balance, i.e. the best combined score.

The log-linear model is more general than the noisy-channel model in that it expresses scoring in terms of aggregation of an unlimited number of feature scores. First a technical note: the log-linear model uses log probabilities rather than regular ones – they are converted to log probabilities simply by applying the log function – and log probabilities are added rather than multiplied.¹ Thus, generally $\log(X \cdot Y) = \log(X) + \log(Y)$ and, more specifically with respect to the noisy-channel model, $\log(P(S|T) \cdot P(T)) = \log(P(S|T)) + \log(P(T))$.

We focus now on the right-hand-side of the log-linear model equation $\sum_{m=1}^M \lambda_m \cdot h_m(T, S)$, where we see that this model comprises a set of log feature scores to be added together. That is, the $\sum_{m=1}^M$ notation indicates that there are a total of M features to be scored and that their individual scores are to be added up (the sigma, \sum). These individual scores are to be computed by multiplying together two feature-specific values, λ_m and $h_m(T, S)$, where λ_m is simply a weight indicating the importance of that feature relative to the other features, and $h_m(T, S)$ is the log probability assigned to the source–candidate pair by that feature. A minimum of two features are usually used: the translation model and language model features exactly as described in (1) for the noisy-channel model. In other words, we can express the noisy-channel model as a log-linear model without changing anything. This is achieved by specifying the set of features to be used as in (4):

m	λ_m	$h_m(T, S)$	
1	0.5	$\log(P(S T))$	(4)
2	0.5	$\log(P(T))$	

Here we have designated feature 1 as the translation model and feature 2 as the language model; the order is not important. The assigned λ values are relative to each other, i.e. the fact that we've given the same values to each of the features indicates that they

are of equal importance² in (4). We then use this feature set to expand the log-linear model formula: as we have two features, there will be a total of two scores to add together and these scores are specified according to the values given in the feature set. Thus, the expansion is as in (5):

$$\begin{aligned} \text{Translation} &= \operatorname{argmax}_T \lambda_1 \cdot h_1(T, S) + \lambda_2 \cdot h_2(T, S) \\ &= \operatorname{argmax}_T 0.5 \log(P(S|T)) + 0.5 \log(P(T)) \end{aligned} \quad (5)$$

A significant advantage of this alternative, log-linear formulation is that we can adjust (or ‘tune’) the relative importance of the features (or ‘parameters’) being used. For example, we could decide that the translation model is twice as important as the language model and should therefore have twice the influence. This change can be made straightforwardly by making the translation model’s λ value double that of the language model: given that the weights must sum to 1, we would give the translation model a weight of 0.67 and the language model a weight of 0.33 (cf. footnote 2). In practice, however, these weights are set empirically relative to the translation task at hand; this is discussed in more detail in Section 3.3.

Another important advantage of the log-linear model is that more features than just these two can be easily incorporated by defining their λ and h values and adding them to the feature set. For example, we could define a translation model which operates at the level of stems rather than full form tokens, or a language model across part-of-speech sequences rather than token sequences. We will not address definitions of additional features any further in this paper, but we provide a list of the typically used features in Section 3.3.

Up to this point, we have been looking at the bigger picture. Now, however, we move to more detailed issues, specifically what the translation and language models actually look like, how they are generated and how they are used during translation. In SMT the resource-induction phase is referred to as *training* while the translation process is referred to as *decoding* (or *testing*). In between, we look at the *tuning* phase where the optimal λ values are learned for the range of features used in an SMT system.

3. Training

The key resources underpinning SMT are corpora: monolingual target-language data for inducing the language model, and sentence-aligned source–target data for inducing the translation model. In SMT, the training data is used to derive the language and translation models, and once these models are generated the parallel data used to induce them are thereafter ignored. When translating new input, the system will access the models but *not* the original training data. In the next sections, we describe how these models are induced.

3.1. LANGUAGE MODEL

As mentioned above, the language model $P(T)$ gives the likelihood that a target string T is actually a valid sentence in the target language. Thus, the language model provides us with two things: (i) a model of the monolingual training corpus and (ii) a method for computing the probability of a (possibly) previously unseen string T using that model. The model extracted from the training corpus comprises relative frequencies for (a subset of) the substrings occurring in that corpus, and the probability assigned to a new sentence is based on the substrings occurring in that sentence which are also part of the model.

One very simple language model comprises just those substrings corresponding to the word tokens in the corpus, where each word type is assigned a probability computed by dividing the number of occurrences of that word in the corpus by the total number of word tokens in the corpus. This model is called a *unigram* model because the substrings extracted comprise single words. An example of a unigram model is given in (6).

$$\begin{array}{ll} \text{Corpus:} & I \text{ need to fly to London tomorrow} \\ \text{Unigram Model:} & P(I)=1/7 \ P(\text{need})=1/7 \ P(\text{to})=2/7 \ P(\text{fly})=1/7 \quad (6) \\ & P(\text{London})=1/7 \ P(\text{tomorrow})=1/7 \end{array}$$

We can use this model to assign a probability to any sentence by simply multiplying together the model's probability for each word in that sentence. We give an example in (7):

$$\begin{aligned} P(I \text{ need to fly tomorrow}) &= P(I) \cdot P(\text{need}) \cdot P(\text{to}) \cdot P(\text{fly}) \cdot P(\text{tomorrow}) \\ &= 1/7 \cdot 1/7 \cdot 2/7 \cdot 1/7 \cdot 1/7 \\ &= 2/16,807 \end{aligned} \quad (7)$$

However, this model has some very significant weaknesses. As the model knows only about single word types, it cannot discriminate between well-formed sentences and those with incorrect word order. For example, the string in (8) is assigned exactly the same probability as the sentence in (7) purely because they contain the same set of lexical items, albeit in a crucially different order.

$$\begin{aligned} P(\text{fly need I to tomorrow}) &= P(\text{fly}) \cdot P(\text{need}) \cdot P(I) \cdot P(\text{to}) \cdot P(\text{tomorrow}) \\ &= 1/7 \cdot 1/7 \cdot 1/7 \cdot 2/7 \cdot 1/7 \\ &= 2/16,807 \end{aligned} \quad (8)$$

For that matter, the model will assign an even higher probability to sentences with greater numbers of frequently occurring words, such as in (9):

$$\begin{aligned} P(\text{to to to to to}) &= P(\text{to}) \cdot P(\text{to}) \cdot P(\text{to}) \cdot P(\text{to}) \cdot P(\text{to}) \\ &= 2/7 \cdot 2/7 \cdot 2/7 \cdot 2/7 \cdot 2/7 \\ &= 32/16,807 \end{aligned} \quad (9)$$

The model also tends to assign higher scores to shorter sentences, simply because the fewer words there are, the fewer probabilities need to be multiplied together. Thus, asked to choose between the proper sentence $P(I \text{ need to fly tomorrow})$, which was assigned probability $2/16,807$ in (7), and the subject-less string $P(\text{need to fly tomorrow})$ as given in (10), the model will prefer the latter.

$$\begin{aligned} P(\text{need to fly tomorrow}) &= P(\text{need}) \cdot P(\text{to}) \cdot P(\text{fly}) \cdot P(\text{tomorrow}) \\ &= 1/7 \cdot 2/7 \cdot 1/7 \cdot 1/7 \\ &= 2/2401 \end{aligned} \quad (10)$$

Of course, if the sentence to be scored contains a word which does not exist in the model then that sentence is assigned probability zero, as in the sentence in (11) which contains the word *return*.

$$\begin{aligned} P(I \text{ need to return tomorrow}) &= P(I) \cdot P(\text{need}) \cdot P(\text{to}) \cdot P(\text{return}) \cdot P(\text{tomorrow}) \\ &= 1/7 \cdot 1/7 \cdot 2/7 \cdot 0 \cdot 1/7 \\ &= 0 \end{aligned} \quad (11)$$

1 However, the model will assign exactly the same score to a sentence containing a single
 2 word which is valid but unseen as it does to all other sentences which contain at least
 3 one unseen word. For example, the (French) sentence in (12) is also assigned a zero
 4 probability.

$$\begin{aligned}
 P(\textit{Je dois retourner demain}) &= P(\textit{Je}) \cdot P(\textit{dois}) \cdot P(\textit{retourner}) \cdot P(\textit{demain}) \\
 &= 0 \cdot 0 \cdot 0 \cdot 0 \\
 &= 0
 \end{aligned}
 \tag{12}$$

9 All of these weaknesses can be addressed, at least to a certain extent. The issue of pre-
 10 ferring shorter sentences is handled in a straightforward way using normalisation (see (13)
 11 and (15) for how bigram and longer n -gram probabilities are normalised by dividing their
 12 counts by the counts of lower order n -grams). The first and most obvious strategy for
 13 dealing with unknown words is to increase the size of the training corpus. However,
 14 despite the fact that language models are commonly induced from corpora comprising
 15 upwards of 1 billion words, of course this strategy alone is not sufficient, as no automati-
 16 cally derived corpus can ever be guaranteed to contain *all* the words in a particular lan-
 17 guage. Techniques referred to as *smoothing* (Jelinek 1977) are therefore also applied in
 18 order to reserve some small amount of probability for unseen words.³ The amount
 19 reserved is likely to be much less than the amount of probability assigned to a word
 20 occurring just once. This means that while no sentence will receive a zero score, the
 21 model can differentiate between sentences containing few (e.g. (11)) and many (e.g. (12))
 22 unseen words.

23 The issue of differentiating between well-formed and ill-formed sentences still remains,
 24 and is addressed by building language models which remember longer substrings from the
 25 training corpus. The *bigram* language model is more complex than the unigram one as it
 26 comprises all two-word substrings (bigrams) occurring in the corpus. Bigram language
 27 models model word *sequences*, meaning that the assigned probabilities reflect the like-
 28 lihood that the second word in the bigram should follow the first word. For instance, the
 29 bigram probability for the word sequence *to fly* will reflect the likelihood that the word
 30 *to* should be followed by the word *fly*, as opposed to being followed by any other word
 31 (*run*, say). Thus, each bigram type is assigned a probability computed by dividing the
 32 number of occurrences of that bigram in the corpus by the total number of occurrences
 33 of the first word of the bigram in the corpus. That is, the bigram probability for *to fly* is
 34 calculated as in (13):

$$\frac{\text{count}(\textit{"to fly"})}{\text{count}(\textit{"to"})}
 \tag{13}$$

38 The bigram model for the same sentence shown in (6) is given in (14). We can use
 39 this model to assign a probability to any sentence by simply multiplying together the
 40 model's probability for each bigram in that sentence.

$$\begin{array}{ll}
 \text{Corpus:} & \textit{I need to fly to London tomorrow} \\
 \text{Bigram Model:} & P(\textit{need|I})=1/1 \ P(\textit{to|need})=1/1 \ P(\textit{fly|to})=1/2 \ P(\textit{to|fly})=1/1 \\
 & P(\textit{London|to})=1/2 \ P(\textit{tomorrow|London})=1/1
 \end{array}
 \tag{14}$$

47 This type of model can be extended to arbitrarily long substrings, or n -grams (of length n).
 48 For example, the trigram language model remembers all three-word substrings in the cor-
 49 pus; the trigram probability for the word sequence *need to fly* will reflect the likelihood

1 that the substring *need to* should be followed by the word *fly*, as opposed to being fol-
 2 lowed by any other word. The trigram probability for *need to fly* is calculated in a similar
 3 manner to (13), as shown in (15):

$$\frac{\text{count}(\text{"need to fly"})}{\text{count}(\text{"need to"})} \quad (15)$$

4
 5
 6
 7 However, we again face the problem of unknowns, this time of unknown n -grams
 8 rather than single words. Of course, unknown n -grams where $n > 1$ are more likely than
 9 unknown words (cf. *to fly* vs. *to above*), so the larger the n -gram language model used,
 10 the worse this problem becomes. For example, the sentence in (16) is assigned zero prob-
 11 ability because the bigram $P(\text{tomorrow}|\text{fly})$ does not exist in the bigram model given in
 12 (14). We saw in (7) that this sentence was assigned a non-zero probability by the unigram
 13 model.

$$\begin{aligned} P(I \text{ need to fly tomorrow}) &= P(\text{need}|\text{I}) \cdot P(\text{to}|\text{need}) \cdot P(\text{fly}|\text{to}) \cdot P(\text{tomorrow}|\text{fly}) \\ &= 1/1 \cdot 1/1 \cdot 1/2 \cdot 0 \\ &= 0 \end{aligned} \quad (16)$$

14
 15
 16
 17
 18 Increasing the size of the training corpus, while still helpful, yields limited returns as
 19 longer n -grams are used, as these are seen less and less often. As before, smoothing tech-
 20 niques are used to assign some very small probability to unseen n -grams, thereby ensuring
 21 that all sentences will obtain non-zero scores. However, this measure does not allow us
 22 to differentiate between unseen n -grams whose individual words occurred in the training
 23 corpus and those containing words never seen before. For example, (16) which contains
 24 the unseen bigram $P(\text{tomorrow}|\text{fly})$ receives exactly the same score as the sentence *I need*
 25 *to fly demain* according to the bigram model in (14), although not according to the uni-
 26 gram model given in (6). Thus, while more context-sensitive models can better reward
 27 correct word order, the less sensitive ones give greater flexibility. In recognition of this,
 28 several models are usually used in combination; that is, for a given sentence, the unigram,
 29 bigram and trigram models might all be used, each with a different weight (typically
 30 preferring larger matches), and their scores combined to give a final probability for that
 31 sentence.

3.2. TRANSLATION MODEL

32
 33
 34
 35 The translation model $P(S|T)$ gives the likelihood that the source sentence S and the can-
 36 didate translation T are translationally equivalent, i.e. that the meaning expressed in S
 37 is also captured in T . Thus, the translation model provides us with two things: (i) a model
 38 of the sentence-aligned source-target training corpus, and (ii) a method for computing
 39 the probability that S and T are equivalent using that model. The model extracted from
 40 the parallel training corpus comprises relative frequencies for (a subset of) the source-
 41 target substring pairs occurring in that corpus. The probability assigned to a new source-
 42 translation pair is then based on the substring pairs occurring in that source-translation
 43 pair which are also part of the model.

44 In fact, the translation model is always *directional*, meaning that $P(S|T)$ actually expresses
 45 the probability that, given target candidate T , source sentence S is likely to have given
 46 rise to T . $P(S|T)$ is, then, a target-to-source (or 'reverse') model of translation. This is
 47 somewhat counter-intuitive: one would expect to want to compute the probability of the
 48 hypothesised target string given the source sentence, i.e. $P(T|S)$. However, the equation
 49 in (1) came about via Bayes theorem to manipulate another equation $P(S, T) = P(T, S)$, 2

\tilde{s}	\tilde{t}	$P(\tilde{s} \tilde{t})$	$P(\tilde{t} \tilde{s})$
the big cat	le grand chat	0.9	0.8
the big cat	le gros chat	0.3	0.2

Fig 1. An example (partial) t-table for English-to-French.

a joint probability, which says that the sentences in the bitext are translations of each other.⁴

The translation model itself comprises a set of target-language phrases, each of which is associated with a list of source-language translations for those phrases and corresponding probabilities. An (invented) example is given in Figure 1. In the first column, we see source-language phrases \tilde{s} ,⁵ with alternative translations \tilde{t} given in the second column. The third and fourth columns give the respective target-to-source and source-to-target probabilities of those phrase pairs. The probabilities are assigned based on relative frequency, i.e. the probability of a source phrase given the target phrase is computed by dividing the number of occurrences of that target phrase in the dataset by the number of occurrences of the aligned source–target phrase pair. Note that with respect to Figure 1, the probabilities in the column headed $P(\tilde{t}|\tilde{s})$ sum to 1, as we assume here that the only translations of ‘the big cat’ are the two French strings under \tilde{t} . By contrast, the probabilities in the column headed $P(\tilde{s}|\tilde{t})$ do *not* sum to 1, as there will be other English translations of ‘le grand chat’ (‘the large cat’, say) and ‘le gros chat’ (e.g. ‘the fat cat’). Note also in regard to this latter that $P(\textit{the fat cat}|\textit{le gros chat})$ is likely to have a much higher probability than $P(\textit{the big cat}|\textit{le gros chat})$, in any representative corpus (‘fat’ being the default dictionary translation of ‘gros’, as opposed to ‘big’).

Before going into detail on how translation models are induced, we first illustrate how they are used to assign probabilities to hypothesised translations. A small snippet of a translation model is given in (17) where we see that each English phrase is associated with a list of French phrases, all of which are associated with probabilities.

I need	→ ...je dois(0.1)...	
I	→ ...je(0.7)...	
need	→ ...dois(0.05)...	
to return	→ ...retourner(0.3)...	(17)
tomorrow	→ ...demain(0.4)...	
return tomorrow by	→ ...retourner demain(0.0001)...	

A translation model probability is then assigned to any source–hypothesis pair by multiplying together the probabilities of the phrase pairs occurring in that pair according to the model. For instance, in (18) three phrase pair probabilities are multiplied together to arrive at a score for the source–hypothesis pair, while in (19) two such probabilities are multiplied together.

Source:	je dois retourner demain	
Hypothesis:	i need to return tomorrow	(18)
$P(\text{Source} \text{Hypothesis})$	$= P(\text{je dois} \text{i need}). P(\text{retourner} \text{to return}).$	
	$P(\text{demain} \text{tomorrow}) = 0.012$	

At this point, the main ‘chicken and egg’ problem for statistical word alignment becomes apparent (Knight 1999b): we have only sentence pairs. If we had access to the probabilistic dictionary and other statistical information mentioned above, we could use this information to work out the most likely word alignments for our sentence pairs. On the other hand, if we had word-aligned sentence pairs we could gather the word correspondence and other statistics we need. In other words, to obtain the statistics we need word-aligned data and to obtain word-aligned data we need the statistics, but we have neither.

Much of the information we need, however, is actually hidden within the sentence-aligned data. Consider, for example, the sentence-aligned pairs in (22). While we have no prior knowledge of the language pair in question, we can see clearly that there is one word common to both source sentences (*ozftsh*) and one word common to both target sentences (*irg*). Thus, it seems reasonable to conclude that these words correspond to each other and should be aligned.⁷

s_1 : gsv ylb	ozftsh	s_2 : z trio	ozftsh	
t_1 : ov tzixlm	irg	t_2 : fmv uroov	irg	(22)

The English–French version of this parallel corpus is given in (23). The common English word is *laughs* and the common French word is *rit*, and it is clear that our deduction that they correspond is correct.

s_1 : the boy	laughs	s_1 : a girl	laughs	
t_1 : le garcon	rit	t_2 : une fille	rit	(23)

This method of deducing word alignments is employed on a large scale in SMT using the Expectation–Maximisation (EM) algorithm (Dempster et al. 1977). That is, the EM algorithm allows us to systematically discover the word alignments for which sufficient evidence exists across the parallel corpus (without recourse to external resources). In the following paragraphs, we explain how EM works using a very simple example.

Consider first the corpus shown in (24(a)) which comprises two sentence pairs. As we currently know nothing about how they should be word-aligned, we simply assume that all possible word alignments are equally likely; the possible word alignments for each sentence pair are given in (24(b)).

(a) Corpus:	green house maison verte	the house la maison			
(b) word alignments:	green house maison verte	green house maison verte	the house la maison	the house la maison	(24)

Before starting the EM algorithm we must first initialise our model as shown in (25). We first collect all the word pairs which exist in the set of word-aligned sentence pairs given in (24(b)). We are assuming that English is the source language and French the target, and that our model gives estimates of $P(\text{source}|\text{target})$ for each word pair identified. Thus, the word pairs extracted from (24(b)) are given down the left side of (25). For instance, we see that *maison* is associated in our dataset with *green*, *house* and *the*. During initialisation, we must also assign some probability to each word pair. As we have no way, as yet, of distinguishing between the suggested English translations for each given French word, we must assume that they are all equally likely. Therefore, as shown in the column labelled ‘0 (Init)’, the three possible translations for *maison* are each assigned prob-

ability $1/3$ while the two possible translations for *verte* are each assigned probability $1/2$, as are the two possible translations for *la*. For any given French word, the probabilities of its set of English translations must always sum to 1.

$$\begin{array}{rcl}
 & & 0 \text{ (Init)} \quad 1 \quad 2 \quad \dots \quad n \\
 P(\text{green}|\text{maison}) & = & \frac{1}{3} \quad \frac{1}{4} \quad \frac{1}{6} \quad \dots \quad 0.01 \\
 P(\text{house}|\text{maison}) & = & \frac{1}{3} \quad \frac{1}{2} \quad \frac{1}{6} \quad \dots \quad 0.98 \\
 P(\text{the}|\text{maison}) & = & \frac{1}{3} \quad \frac{1}{4} \quad \frac{1}{6} \quad \dots \quad 0.01 \\
 P(\text{house}|\text{verte}) & = & \frac{1}{2} \quad \frac{1}{2} \quad \frac{1}{3} \quad \dots \quad 0.01 \\
 P(\text{green}|\text{verte}) & = & \frac{1}{2} \quad \frac{1}{2} \quad \frac{1}{3} \quad \dots \quad 0.99 \\
 P(\text{house}|\text{la}) & = & \frac{1}{2} \quad \frac{1}{2} \quad \frac{1}{3} \quad \dots \quad 0.01 \\
 P(\text{the}|\text{la}) & = & \frac{1}{2} \quad \frac{1}{2} \quad \frac{1}{3} \quad \dots \quad 0.99
 \end{array} \tag{25}$$

At this point, the initialisation procedure is complete. The EM algorithm is an iterative algorithm, meaning that it will be carried out a number of times, and the probabilities given in the model are updated after every iteration. The remaining columns in (25) correspond to successive iterations of the algorithm, i.e. column 1 gives the probabilities as they are at the end of the first iteration, column 2 gives them as they are at the end of the second iteration and column n gives the probabilities at the end of the final iteration.

Each iteration comprises two steps which can be referred to as the E(xpectation)-step and the M(aximisation)-step. During the E-step, we assign a probability to the alternative word alignments for each of the sentence pairs in the corpus using the word pair probabilities specified in the model. During the M-step, we compute new probabilities for the word pairs in the model based on the probabilities of the word alignments specified in the corpus, and update the model using these new word pair probabilities. We will further explain these steps by showing how the probabilities in column 1 in (25) are computed given the word alignments in (24(b)) and the word pair probabilities in column 0 in (25); the steps are illustrated in (26).

First the E-step: in order to compute the probability of a particular word alignment for a sentence pair, we multiply together the probabilities (given in the model) of each linked source word given the target word to which it is linked. Thus, for example, for the first word-aligned sentence pair shown in (26), we see that the probability of *green* given *maison* is $1/3$ and the probability of *house* given *verte* is $1/2$ according to column 0 of our model. Multiplied together, they give a probability of $1/6$ for the overall word alignment. Similarly, for the other possible word alignment of that same sentence pair, we see that the probability of *house* given *maison* is $1/3$ and the probability of *green* given *verte* is $1/2$, thus also giving an overall probability of $1/6$ for the word alignment. Once this calculation has been completed for all the word alignments in the dataset, a final piece of book-keeping remains: wherever we have more than one possible word alignment for a sentence pair, we must normalise such that the probabilities of those alternative word alignments sum to 1. To achieve this, we add up the probabilities of the alternative alignments for the given sentence pair and then divide each alternative word alignment by the total. For instance, as we have two different word alignments for the sentence pair $\langle \text{green house, maison verte} \rangle$ and each has probability $1/6$, we divide both those probabilities by $2/6$ to give the final probability for each. Not surprisingly, given the probabilities in column 0 of our model, both these alignments are judged to be equally likely at this point (each having probability of $1/2$). However, our first iteration is still not complete as we have not yet addressed the M-step.

<u>Iteration 1 – E-step:</u>	<u>Iteration 1 – M-step:</u>
$P \left(\begin{array}{c} \text{green house} \\ \quad \\ \text{maison verte} \end{array} \right) = \frac{1}{3} \cdot \frac{1}{2} = \frac{1}{6} \quad / \frac{2}{6} = \frac{1}{2}$	$P(\text{green} \text{maison}) = \frac{1}{2} \quad / 2 = \frac{1}{4}$
$P \left(\begin{array}{c} \text{green house} \\ \times \\ \text{maison verte} \end{array} \right) = \frac{1}{3} \cdot \frac{1}{2} = \frac{1}{6} \quad / \frac{2}{6} = \frac{1}{2}$	$P(\text{house} \text{maison}) = \frac{1}{2} + \frac{1}{2} = 1 \quad / 2 = \frac{1}{2}$
$P \left(\begin{array}{c} \text{the house} \\ \quad \\ \text{la maison} \end{array} \right) = \frac{1}{2} \cdot \frac{1}{3} = \frac{1}{6} \quad / \frac{2}{6} = \frac{1}{2}$	$P(\text{the} \text{maison}) = \frac{1}{2} \quad / 2 = \frac{1}{4}$
$P \left(\begin{array}{c} \text{the house} \\ \times \\ \text{la maison} \end{array} \right) = \frac{1}{2} \cdot \frac{1}{3} = \frac{1}{6} \quad / \frac{2}{6} = \frac{1}{2}$	$P(\text{house} \text{verte}) = \frac{1}{2} \quad / 1 = \frac{1}{2}$
	$P(\text{green} \text{verte}) = \frac{1}{2} \quad / 1 = \frac{1}{2}$
	$P(\text{house} \text{la}) = \frac{1}{2} \quad / 1 = \frac{1}{2}$
	$P(\text{the} \text{la}) = \frac{1}{2} \quad / 1 = \frac{1}{2}$

During the M-step, we re-estimate the model's probabilities given the relative frequency of occurrence of each word pair in the set of word alignments. Crucially, however, we do not count a frequency of 1 each time we see a link in the dataset between the word pair we are interested in. Instead, we count the probability that we associated with that particular word alignment during the E-step. This is illustrated in the M-step part of (26). Consider, for example, the word pair *green* and *maison*. This word pair is linked in a single-word alignment in the dataset (the first one), and this word alignment has probability 1/2 according to the E-step. Thus, the frequency of *green* given *maison* is 1/2. We can contrast this with the word pair *house* and *maison*: this word pair is linked in two different word alignments in the dataset, each of which has probability 1/2. Thus, we add together the probabilities of those word alignments, meaning that the frequency of *house* given *maison* is 1. At this point, we still have frequencies rather than probabilities. In order to obtain probabilities, we add up the frequencies of the alternative translations for each target word and then divide each alternative translation frequency by the total. For instance, as we have three different translations for the word *maison* and their frequencies sum to 2, we divide each frequency by 2 to give a probability. Once this has been done, the M-step is complete. The model probabilities in (25) are updated by simply replacing the existing probabilities with those output from the M-step – as shown in column 1 – and iteration 1 is complete.

After just one iteration, we can clearly see that we have improved our knowledge regarding how likely particular word pairs are given the dataset. The three proposed translations for *maison* – *green*, *house* and *the* – are no longer equally likely. Rather, our model now says that, given *maison*, *house* is twice as likely as either *green* or *the*. However, our model still tells us that *house* and *green* are still equally likely given *verte*, and that *house* and *the* are still equally likely given *la*. Thus, we repeat the process in the hope that further improvements in our estimates can be made.

The E- and M-steps for iteration 2 are shown in (27). During the E-step, the probabilities of the different word alignments are computed in exactly the same way as before. This time, however, when we multiply together the probabilities of each linked source word given the target word to which it is linked, those probabilities come from the *updated* model, i.e. from column 1 in (25) rather than column 0. Thanks to these updated probabilities, we can now differentiate between the different word alignments for each sentence pair: the second word alignment is preferred over the first and the third is preferred over the fourth.

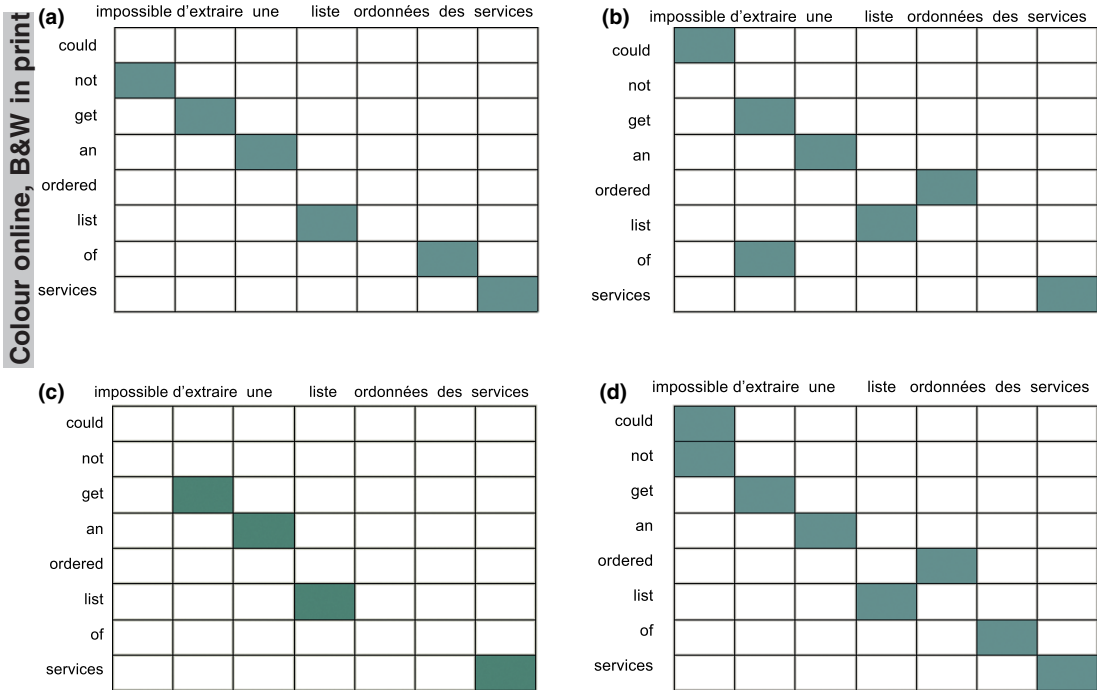
Iteration 2 – E-step:	Iteration 2 – M-step:
$P \begin{pmatrix} \text{green house} \\ \quad \\ \text{maison verte} \end{pmatrix} = \frac{1}{4} \cdot \frac{1}{2} = \frac{1}{8} \quad / \frac{3}{8} = \frac{1}{3}$	$P(\text{green} \text{maison}) = \frac{1}{3} \quad / 2 = \frac{1}{6}$
$P \begin{pmatrix} \text{green house} \\ \diagdown \quad \diagup \\ \text{maison verte} \end{pmatrix} = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4} \quad / \frac{3}{8} = \frac{2}{3}$	$P(\text{house} \text{maison}) = \frac{2}{3} + \frac{2}{3} = \frac{4}{3} \quad / 2 = \frac{4}{6}$
$P \begin{pmatrix} \text{the house} \\ \quad \\ \text{la maison} \end{pmatrix} = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4} \quad / \frac{3}{8} = \frac{2}{3}$	$P(\text{the} \text{maison}) = \frac{1}{3} \quad / 2 = \frac{1}{6} \quad (27)$
$P \begin{pmatrix} \text{the house} \\ \diagdown \quad \diagup \\ \text{la maison} \end{pmatrix} = \frac{1}{2} \cdot \frac{1}{4} = \frac{1}{8} \quad / \frac{3}{8} = \frac{1}{3}$	$P(\text{house} \text{verte}) = \frac{1}{3} \quad / 1 = \frac{1}{3}$
	$P(\text{green} \text{verte}) = \frac{2}{3} \quad / 1 = \frac{2}{3}$
	$P(\text{house} \text{la}) = \frac{1}{3} \quad / 1 = \frac{1}{3}$
	$P(\text{the} \text{la}) = \frac{2}{3} \quad / 1 = \frac{2}{3}$

These probabilities in turn feed into the re-estimation of the word pair probabilities in the M-step. For example, although the word pairs *house* given *verte* and *green* given *verte* each occur in a single-word alignment, it is the word alignment probabilities which are counted. Therefore, the latter is counted 2/3 times whereas *house* given *verte* is counted just 1/3 times. After the M-step, the model is updated, allowing us to complete column 2 in (25). Again, we have improved our knowledge: the proposed translations for *verte* are no longer equally likely, but rather *green* is twice as likely as *house*. We continue iterating until we stop improving our estimates; column *n* in (25) suggests the kind of model we may end up with when the process is complete.

3.2.2. Phrase Alignment

Two important characteristics of the word alignment process described above are as follows. Firstly, the alignment algorithm is directional, i.e. if we have an English–French parallel dataset, word-aligning it with English as source and French as target will not necessarily yield exactly the same result as word-aligning it with French as source and English as target. Secondly, the alignment algorithm can produce 1-to-1 alignments (i.e. with one source word mapping to one target word), and 1-to-many alignments (one source word mapping to more than one target word) but not many-to-1 alignments and many-to-many alignments. In practice when training an SMT system, the parallel corpus is word-aligned in both directions. This means that even when the objective is to build a French-to-English system, both French-to-English and English-to-French alignments are generated. This allows us to also extract many-to-1 alignments from the target-to-source alignment. However, we still do not have many-to-many alignments, i.e. source-target phrase pairs.

Phrase-alignment heuristics (Koehn et al. 2003; Och and Ney 2003) which operate on the output of the word alignment algorithm are frequently used to extract these many-to-many alignments. Firstly, word alignment is performed in both directions on each training sentence pair and mapped to a bitext grid as shown in Figures 2a,b. We then compute a more precise set of word alignments by taking only those that occur in both sets, as shown in Figure 2c. From this point, it is possible to proceed by iteratively adding further alignments where (i) the alignments are adjacent, and (ii) either the source or target word is currently unaligned: this is shown in Figure 2(d). Finally, phrase pairs can be read from this extended grid: examples given in Groves and Way (2005) include the following:



18 Fig 2. Extracting phrase alignments from word alignments (Groves and Way 2005). The example was taken from a parallel corpus, hence the grammatical error in the French string. (a) English–French word alignment, (b) French–English word alignment, (c) Intersection of alignments and (d) Intersection extended to union.

- could not ↔ impossible
- could not get ↔ impossible d'extraire
- get an ↔ d'extraire une
- ordered list ↔ liste ordonnées
- get an ordered list ↔ d'extraire une liste ordonnées
- could not get an ordered list ↔ impossible d'extraire une liste ordonnées.

3.3. TRANSLATION AND EVALUATION FOR TRAINING PURPOSES

Once we have learned the required models from the training data, and assuming we have an algorithm that accepts an input sentence and returns a translation (to be discussed in detail in Section 4), a final training task still remains before we have a full system that can then be used to translate. In Section 2, we looked at the righthand-side of the log-linear model equation $\sum_{m=1}^M \lambda_m \cdot h_m(T, S)$ and we saw that it comprises a set of log feature scores to be added together. Recall that the feature scores are computed by multiplying together two feature-specific values, λ_m and $h_m(T, S)$, where λ_m is a weight indicating the importance of that feature relative to the other features being used. By changing the λ values, we can adjust (or 'tune') the relative importance of the features being used. While we have focused on the language model and translation model features in this paper, other features are frequently used in addition to these, typically including:

1. an n -gram language model over target sequences,
2. a source-to-target phrase table,

3. a target-to-source phrase table,
4. source-to-target lexical translation probabilities,
5. target-to-source lexical translation probabilities,
6. a phrase reordering model,
7. the standard word/phrase penalty which allows for control over the length of the target sentence.

Each of these features must be associated with a λ value controlling its influence on the final translation decision. Och (2003) proposed a widely adopted approximation technique called Minimum Error Rate Training (MERT) to estimate the model parameters for a small number of features, such as those listed above.

MERT requires that we define an error function that corresponds to translation accuracy. Most commonly, the error function used is the BLEU evaluation metric (Papineni et al. 2001, 2002), which works by comparing each output translation to a reference translation and computing a similarity/dissimilarity score; we describe BLEU in detail in Section 3.3.1. MERT also requires a second component – a development set (or ‘devset’) – which is a set of source–target sentence pairs that were not part of the training set used to induce the translation model (or any other features), and that are not part of any test set that might be used to evaluate overall system performance. For any set of λ values, we translate the source sentences in the devset via the MT system and use BLEU to compare the output translations against the devset target sentences. This gives us a measure of translation accuracy, where we assume the target sentences in the devset to be the correct ‘reference’ translations that we are aiming to produce from the MT system. Using MERT, if we then adjust the λ values, we can again translate and evaluate using the devset to see whether the BLEU score has gone up or down. Thus, MERT is an iterative process: we repeatedly translate the devset using different λ values in order to identify the set of λ values that gives the best performance according to our error function on that devset.

Essentially, therefore, MERT involves optimising system performance to *one* particular metric on a devset of sentences, and hoping that this carries forward to the test set at hand; the closer the devset to the test set, the more appropriate the settings obtained via MERT will be. MERT provides a simple and efficient method to estimate the optimal weights for the features in the SMT model; however, it can only handle a small number of parameters, and when the number of features increases, there is no guarantee that MERT will find the most suitable combination (Chiang et al. 2008).

3.3.1. The BLEU Metric

In this short section, we digress into the related field of MT evaluation in order to describe how the evaluation metric most commonly used for MERT works.

The BLEU metric (Papineni et al. 2001, 2002) evaluates MT system quality by comparing output translations to their reference translations in terms of the numbers of co-occurring n -grams. The main score calculated is the n -gram precision p_n for each pair of candidate and reference sentences. This score represents the proportion of n -word sequences in the candidate translation which also occur in the reference translation. Importantly, if an n -gram occurs j times in the candidate translation and i times in the reference translation such that $i \leq j$ then this sequence is counted only i times; this corresponds to the intuition that ‘a reference word sequence should be considered exhausted after a matching candidate word sequence has been identified’ (Papineni et al. 2001). Thus, n -gram precision p_n is calculated according to equation (28):

$$p_n = \frac{|c_n \cap r_n|}{|c_n|} \quad (28)$$

where

- c_n is the multiset of n -grams occurring in the candidate translation.
- r_n is the multiset of n -grams occurring in the reference translation.
- $|c_n|$ is the number of n -grams occurring in the candidate translation.
- $|c_n \cap r_n|$ is the number of n -grams occurring in c_n that also occur in r_n such that elements occurring j times in c_n and i times in r_n occur maximally i times in $|c_n \cap r_n|$.

As it is generally not the case that MT output is evaluated one sentence at a time, n -gram precision is usually calculated over large sets of sentences (typically not less than 1000). In this case, p_n is the proportion of co-occurring n -word sequences in the set over the total number of n -word sequences in the set.

While precision scores p_n can be obtained for any value of n ,⁸ Papineni et al. (2001) point out that greater robustness can be achieved by combining scores for all values of n into a single metric. It is not surprising that as the value of n increases, the score p_n decreases because longer matching word sequences are more difficult to find (cf. Section 3.1 for a similar discussion regarding higher-order n -grams in language modelling). If the average n -gram precision score is calculated without taking this factor into account (i.e. by simply summing the values for p_n and dividing by N , the largest value for n), then the scores for longer n -grams will be too small to have much influence on the final score. In order to make the BLEU metric more sensitive to longer n -grams, the combined score p_N is calculated by summing over the logarithm of each p_n multiplied by weight $1/N$ as given in equation (28):

$$p_N = \exp\left(\sum_{n=1}^N \frac{1}{N} \log(p_n)\right) \quad (29)$$

A candidate translation which is longer than its reference translation is implicitly penalised during the calculation of p_n . In order to impose a corresponding penalty on candidate translations which are shorter than their reference translations, a *brevity penalty* BP is introduced and the combined precision score p_N is multiplied by this penalty. BP is incorporated to prevent the output of single-word utterances like 'the', which would otherwise score highly. Papineni et al. (2001) state that BP is a decaying exponential in the length of the reference sentence over the length of the candidate sentence. This means that if the reference is the same length or longer than the candidate, then the penalty is 1, and greater than 1 if the candidate is shorter than the reference. Furthermore, if candidate c_x is 1 word shorter than its reference r_x and c_y is also 1 word shorter than r_y , but r_x is longer than r_y , then the BP for c_y should be greater than the BP for c_x . Thus, BP is calculated according to equation (30):

$$BP = e^{\max\left(1 - \frac{\text{length}(R)}{\text{length}(C)}, 0\right)} \quad (30)$$

Note that as calculating the brevity penalty for each sentence and averaging it over the set of sentences is considered by Papineni et al. (2001) to be unduly harsh, it is computed over the entire corpus, i.e. $\text{length}(R)$ is the number of words in the reference set and $\text{length}(C)$ the number of words in the candidate set. This penalty is then applied to the precision score for the entire candidate translation corpus according to equation (31):

$$BLEU = BP \cdot p_N \quad (31)$$

As the ranking behaviour is more visible in the log domain, Papineni et al. (2001) give equation (32):

$$\log BLEU = \min\left(1 - \frac{\text{length}(R)}{\text{length}(C)}, 0\right) + \sum_{n=1}^N \frac{1}{N} \log(p_n) \quad (32)$$

In summary, the combination of n -gram precision and penalties for shorter translations mean that in order to achieve a high BLEU score, a set of candidate translations must match the reference translations in length, in word choice and in word order Papineni et al. (2001).

Other evaluation metrics based on similar principles as BLEU have also been introduced, e.g. NIST (Doddington 2002), followed soon thereafter by GTM (Turian et al. 2003), and a couple of years later by METEOR (Banerjee and Lavie 2005). Despite a number of researchers pointing out various problems with such metrics (cf. Callison-Burch et al. (2006) for a criticism of BLEU), all these metrics are still actively used in MT research. While a wide variety of techniques, including human evaluation, are used to compare systems and assess overall translation quality, automatic similarity-based metrics are critical in facilitating the use of MERT.

4. Translation

In previous sections we showed how SMT language and translation models are induced from data (training), and how these and any other parameters in the log-linear model can be optimised (tuning). In Section 3.3, we explicitly assumed that we have an algorithm that accepts an input sentence and returns the most likely translation, but we did not provide any details. In this section, we outline this final piece of the puzzle: the translation phase, referred to also as ‘decoding’ (or ‘testing’).

Given any source sentence and a hypothesised translation, we can compute a corresponding probability by using the models generated during training in relation to equation (2) for the log-linear model. By extension, given any source sentence and a set of hypothesised translations, we can work out which hypothesis is most likely. Decoding, therefore, is a search problem: we must find the most probable translation amongst all translations possible given the model.

In theory, this would involve generating the full set of possible translations for any input string. Translations are generated by simply matching substrings from the input sentence against the translation model and, where available, retrieving their translations. These substring translations can then be concatenated in any order to generate a full translation hypothesis. Figure 3 shows some of the possible substring translations for a Spanish sentence, from which any sequence of substrings which fully covers the source sentence can be selected in any order to form a hypothesis.

There is no necessity in the decoding process to start with the leftmost source word, and proceed on a strictly left-to-right basis. In our Spanish-to-English example it would seem like a reasonable strategy but of course for some language pairs (Arabic-to-English, say) that would be entirely the wrong thing to do. The good thing for the decoder is that it essentially makes no difference, and the best translation – based on the probability computed using the log-linear formula – will be chosen irrespective of which source-language words are translated in which order.

Maria	no	daba	una	bofetada	a	la	bruja	verde
Mary	not	give	a	slap	to	the	witch	green
	did not		a	slap	by		green	witch
	no		slap		to the			
	did not give				to			
					the			
			slap			the	witch	

Fig 3. Some English translation options for the Spanish sentence *Maria no daba una bofetada a la bruja verde* (Koehn 2004).

However, it was shown in Knight (1999a) that the decoding problem is NP-complete, i.e. that it cannot be guaranteed that we can generate and score all hypotheses for a source sentence in a finite amount of computing time. Accordingly, decoding is also an optimisation problem; we must avoid generating target-language hypotheses which are unlikely to end up being the most probable, without also missing the good solutions by accident. In short, we must try to find the most probable translation without actually generating all of the hypotheses.

The state-of-the-art implementation of decoding for SMT is a beam-search decoder (Koehn et al., 2003, 2007). A ‘beam’ is just an arbitrary number of hypotheses (say, 10,000) that are maintained at any point in the translation process, and such a cut-off ensures that the runtime of any system is manageable in practice. At the start of the translation process, no words have yet been translated. Source-language words are expanded in a monotone or non-monotone manner, i.e. following the source word/phrase order or not. New hypotheses are generated from the existing expanded hypotheses by extending them with a phrasal translation that covers some of the source input words which have not yet been translated. Each hypothesis is added into a beam stack as a new node. Figure 4 shows six stacks, which is sufficient for covering the first 6 words in target-language hypotheses. Each stack is marked by the covered source words (one word covered in stack 1, two words in stack 2, etc.) during expansion. A newly created hypothesis will be placed in a new stack further down. For example, if we return to Figure 3, the second from the top phrase in stack 2 (comprising two words, *Mary did*, say) is linked to various hypotheses in stacks 3 (*not*, i.e. three words – *Mary did not* – are now covered), 4 (*give*, four words: *Mary did not give*) and 5 (*a*, five words: *Mary did not give a*), all of which might be partial candidate translations (i.e. valid paths through the search space in Figure 3) of the source string *Maria no daba una bofetada a la bruja verde*.

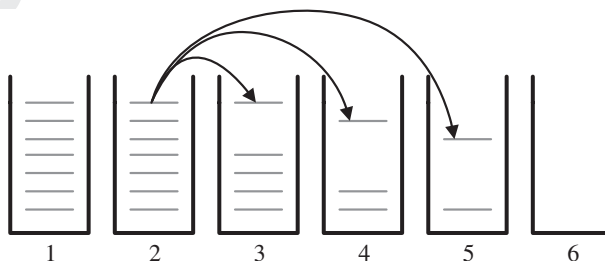


Fig 4. Hypothesis expansion via stack decoding (Koehn 2004).

1 The key idea is that each beam stack has a maximum size, and once the maximum size
 2 is reached a new hypothesis can be added only if its score is higher than the lowest-scoring
 3 hypothesis on the stack (in which case that lowest-scoring hypothesis is dropped to
 4 make room). Of course, hypotheses are scored in part by computing how likely they are
 5 according to the log-linear formula. However, relying on just this probability is not suffi-
 6 cient for pruning purposes. Stack 4, for instance, holds the hypothesised four-word trans-
 7 lation sequences from the source sentence. Some of those word sequences will have been
 8 easier to translate than others, meaning they will likely have higher probabilities. How-
 9 ever, we do not want to prune away hypotheses from the stack that have lower probabili-
 10 ties simply because they are more difficult to translate, leaving only translations for the
 11 easier sequences behind. Therefore, when scoring a hypothesis to decide its inclusion on
 12 a stack we also factor in the *future translation cost* (Koehn 2010:167–72); this reflects how
 13 difficult it might be to translate the source words *not* covered by that hypothesis. Thus,
 14 the future cost estimation should favour hypotheses that have already covered difficult
 15 parts of the sentence and have only easy parts left, while discounting hypotheses that have
 16 covered the easy parts first.⁹

17 The final states in the search are hypotheses that cover all source words: these are
 18 found in stack n where n is the length of the source sentence (9, in Figure 3). Among
 19 the n -best list of possible translations created by the decoder following different paths
 20 through the search space, the one with the highest probability is finally selected as the
 21 best translation.

22 5. Conclusion

23 In this paper, we have given an overview of the basic model of SMT, along with in-
 24 depth presentations of the core components of SMT systems. In Section 2 we presented
 25 the big-picture view of the SMT approach. In Section 3 we described the main
 26 approaches to inducing training data from corpora, focusing in particular on language
 27 modelling, translation modelling and tuning via Minimum Error Rate Training. Finally,
 28 in Section 4 we gave an overview of the decoding process used in SMT systems to gen-
 29 erate translations for previously unseen text. For a detailed discussion on the role of lin-
 30 guists, translators and translations in SMT, along with a description of how the ideas
 31 behind SMT developed and how they relate to other MT approaches, please see Way
 32 and Hearne (2010). For those readers who are interested in more in-depth accounts of
 33 SMT, we recommend Brown et al. (1993), Knight (1999b) and Knight (2010). For **3**
 34 pointers to the primary literature relating to each of the sub-tasks in SMT, as well as
 35 other MT paradigms, consult section 7 of Way (2010), or the appropriate sections in
 36 Koehn (2010).

37 Short Biographies

38 ????????

39 Acknowledgement

40 This work was partially funded by a number of Science Foundation Ireland (<<http://www.sfi.ie>>)
 41 awards, namely: Principal Investigator Award 05/IN/1732, Basic Research **4**
 42 Award 06/RF/CMS064 and CSET Award 07/CE/I1142. Thanks also to Philipp Koehn
 43 for making available some of his original figures.

Notes

* Correspondence address: Andy Way, NCLT & CNGL, School of Computing, Dublin City University, Collins Avenue, Glasnevin, Dublin D9, Ireland. E-mail: away@computing.dcu.ie

¹ This is done for a very simple reason: the number of multiplications of probabilities performed is often very large. In almost all cases, these probabilities are less than 1, so the resultant score is often extremely small, and cannot be accurately captured unless a large number of decimal places are used.

² Note that these λ values sum to 1. In practice, as long as two equivalently weighted features have the same value, and that normalisation (cf. (13) and (15) below) later takes place, this requirement that the λ values sum to 1 can be relaxed.³ The main differences between different smoothing techniques relate to how much probability mass is subtracted out ('discounting') and how it is redistributed ('back-off'). The most popular method used in SMT is Kneser–Ney smoothing (Kneser and Ney 1995).

⁴ In Section 3.3, the reader will note that both a source-to-target and a target-to-source translation model (also known as a *phrase table*, *translation table* or indeed *t-table*) are employed as typical features used in an SMT system (cf. Figure 1). While there is little justification *per se* for including both – translation is, after all, directional by its very nature – doing so causes translation quality to improve. Anecdotally, the story goes that one of the major SMT developers in our field one day selected the wrong translation table, but that the system's BLEU score didn't go down by very much. Once the error was discovered, including both t-tables caused performance to improve upon using just the one, correct t-table, and henceforth we all use both!

⁵ The 'n' indicates that we are dealing with 'phrases', as opposed to 'sentences'; note that the term 'phrase' is not used in any linguistic sense here, but rather is used to refer to any substring (*n*-gram sequence) of between 1 and (usually around) 7 words in length.

⁶ <<http://www.fjoch.com/GIZA++.html>>.

⁷ Our exposition here draws heavily on how one might translate between Arcturan and Centauri (Knight 1997, 1999b). For newcomers to SMT, these resources are an excellent starting point.

⁸ Scores can be obtained for any *reasonable* value of *n*; in Papineni et al. (2001, 2002) the maximum value for *n* considered was 4.

⁹ The 'ease' or 'difficulty' associated with translating certain parts of a sentence is usually expressed in terms of weighted log probabilities which take into account (at least) language model, translation model and reordering costs. As you might expect, common words are 'easier' to translate in this model than less frequent words, despite these being among the 'hardest' words to get right for humans.

Works Cited

- Banerjee, S. and A. Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization, 65–72. Ann Arbor, MI. **6**
- Brown, P., J. Cocke, S. Della Pietra, V. Della Pietra, F. Jelinek, J. Lafferty, R. Mercer, and P. Roosin. 1990. A statistical approach to machine translation. Computational Linguistics 16. 79–85.
- , S. Della Pietra, V. Della Pietra, and R. Mercer. 1993. The mathematics of statistical machine translation: parameter estimation. Computational Linguistics 19. 263–311.
- Callison-Burch, C., M. Osborne, and P. Koehn. 2006. Re-evaluating the role of BLEU in machine translation research. EACL-2006 11th Conference of the European Chapter of the Association for Computational Linguistics, Proceedings. 249–56. Trento, Italy. **7**
- Chiang, D., Y. Marton, and P. Resnik. 2008. Online large-margin training of syntactic and structural translation features. Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing, 224–33. Honolulu, HI.
- Dempster, A., N. Laird, and D. Rubin 1977. Maximum likelihood from incomplete data via the EM algorithm. Journal of the Royal Statistical Society Series B 39. 1–38.
- Doddington, G. 2002. Automatic evaluation of MT quality using N-gram co-occurrence statistics. Proceedings of Human Language Technology Conference 2002, 138–45. San Diego, CA. **8**
- Groves, D. and A. Way 2005. Hybrid data-driven models of machine translation. Machine Translation 19. 301–23.
- Jelinek, F. 1977. Statistical methods for speech recognition. Cambridge, MA: MIT Press.
- Kneser, R. and H. Ney. 1995. Improved backing-off for m-gram language modeling. Proceedings of the IEEE International Conference on Acoustics Speech and Signal Processing, Vol. 1, 181–4, Detroit, MI.
- Knight, K. 1997. Automating knowledge acquisition for machine translation. AI Magazine 18. 81–96.
- , 1999a. Decoding complexity in word-replacement translation models. Computational Linguistics 25. 607–15.
- , 1999b. A statistical machine translation tutorial workbook. <<http://www.isi.edu/natural-language/mt/wkbk.rtf>>.

- 1 Koehn, P. 2004. 'Pharaoh.' A beam search decoder for phrase-based statistical machine translation models. Machine
 2 Translation: from Real Users to Research, 6th Conference of the Association for Machine Translation in the
 3 Americas, AMTA 2004, Proceedings, 115–24, Washington, DC. **9**
- 4 ——. 2010. Statistical machine translation. Cambridge, UK: Cambridge University Press.
- 5 ——, F. Och, and D. Marcu. 2003. Statistical phrase-based translation. HLT-NAACL: Human Language Technol-
 6 ogy Conference of the North American Chapter of the Association for Computational Linguistics, 127–33.
 7 Edmonton, AL, Canada. **10**
- 8 ——, H. Hoang, A. Birch, C. Callison-Burch, M. Federico, N. Bertoldi, B. Cowan, W. Shen, C. Moran,
 9 R. Zens, C. Dyer, O. Bojar, A. Constantin, and E. Herbst. 2007. Moses: open source toolkit for statistical
 10 machine translation. Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics
 11 Companion Volume Proceedings of the Demo and Poster Sessions, 177–80. Prague, Czech Republic. **11**
- 12 Och, F. 2003. Minimum error rate training in statistical machine translation. 41st Annual Meeting of the Associa-
 13 tion for Computational Linguistics, 160–7. Sapporo, Japan. **12**
- 14 ——, and H. Ney. 2002. Discriminative training and maximum entropy models for statistical machine translation.
 15 40th Annual Meeting of the Association for Computational Linguistics, 295–302. Philadelphia, PA. **13**
- 16 ——, and ——. 2003. A systematic comparison of various statistical alignment models. Computational Linguistics
 17 29. 19–51.
- 18 Papineni, K., S. Roukos, T. Ward, and W.-J. Zhu. 2001. BLUE: a method for automatic evaluation of machine
 19 translation. Technical report, IBM T.J. Watson Research Center, Yorktown Heights, NY.
- 20 ——, ——, ——, and ——. 2002. BLEU: a method for automatic evaluation of machine translation. 40th Annual
 21 Meeting of the Association for Computational Linguistics, 311–8. Philadelphia, PA. **14**
- 22 Turian, J., L. Shen, and D. Melamed. 2003. Evaluation of machine translation and its evaluation. Machine transla-
 23 tion summit IX, ed. by ???, 386–93. New Orleans, LA: ????. **15**
- 24 Way, A. 2010. Machine translation. Handbook for computational linguistics and natural language processing, ed. by
 25 S. Lappin, C. Fox, and A. Clark. Chichester, UK: Wiley Blackwell. **16**
- 26 ——, and M. Hearne. forthcoming. On the role of translations in state-of-the-art statistical machine translations.
 27 *Compass*. **17**

Author Query Form

Journal: LNC3

Article: 274

Dear Author,

During the copy-editing of your paper, the following queries arose. Please respond to these by marking up your proofs with the necessary changes/additions. Please write your answers on the query sheet if there is insufficient space on the page proofs. Please write clearly and follow the conventions shown on the attached corrections sheet. If returning the proof by fax do not write too close to the paper's edge. Please remember that illegible mark-ups may delay publication.

Many thanks for your assistance.

Query reference	Query	Remarks
Q1	AUTHOR:A running head short title was not supplied; please check if this one is suitable and, if not, please supply a short title of up to 40 characters that can be used instead.	
Q2	AUTHOR:Can 'equation in (1)' be changed to 'equation (1)'? Please confirm.	
Q3	AUTHOR:Knight (2010) has not been included in the list, please supply publication details.	
Q4	AUTHOR:Please check this website address and confirm that it is correct.	
Q5	AUTHOR:Please provide short biographies for both the authors as per the journal style.	
Q6	AUTHOR:Please provide the publisher for reference 'Banerjee and Lavie (2005)'.	
Q7	AUTHOR:Please provide the publisher for reference 'Callison-Burch 2006'.	
Q8	AUTHOR:Please provide the publisher for reference 'Doddington 2002'.	
Q9	AUTHOR:Please provide the publisher for reference 'Koehn 2004'.	
Q10	AUTHOR:Please provide the publisher for reference 'Koehn et al. 2003'.	

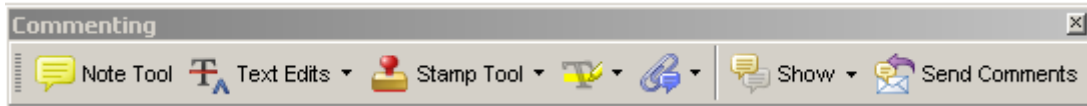
Q11	AUTHOR: Please provide the publisher for reference 'Koehn et al. 2007'.	
Q12	AUTHOR: Please provide the publisher for reference 'Och 2003'.	
Q13	AUTHOR: Please provide the publisher for reference 'Och and Ney 2002'.	
Q14	AUTHOR: Please provide the publisher for reference 'Papineni et al. 2002'.	
Q15	AUTHOR: Please provide the Editors and name of the publisher for this book.	
Q16	AUTHOR: Please provide page range for reference 'Way 2010'.	
Q17	AUTHOR: Please provide the correct journal name, volume number and page range for reference 'Way and Hearne 2010'.	
Q18	AUTHOR: Figure 2 has been saved at a low resolution of 169 dpi. Please resupply at 600 dpi. Check required artwork specifications at http://authorservices.wiley.com/submit_illustration.asp?site=1	

USING E-ANNOTATION TOOLS FOR ELECTRONIC PROOF CORRECTION

Required Software

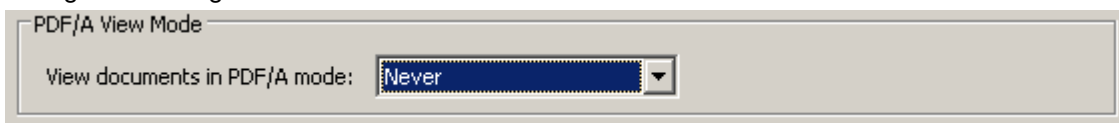
Adobe Acrobat Professional or Acrobat Reader (version 7.0 or above) is required to e-annotate PDFs. Acrobat 8 Reader is a free download: <http://www.adobe.com/products/acrobat/readstep2.html>

Once you have Acrobat Reader 8 on your PC and open the proof, you will see the Commenting Toolbar (if it does not appear automatically go to Tools>Commenting>Commenting Toolbar). The Commenting Toolbar looks like this:



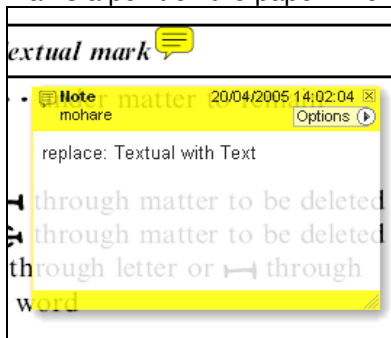
If you experience problems annotating files in Adobe Acrobat Reader 9 then you may need to change a preference setting in order to edit.

In the “Documents” category under “Edit – Preferences”, please select the category ‘Documents’ and change the setting “PDF/A mode:” to “Never”.



Note Tool — For making notes at specific points in the text

Marks a point on the paper where a note or question needs to be addressed.

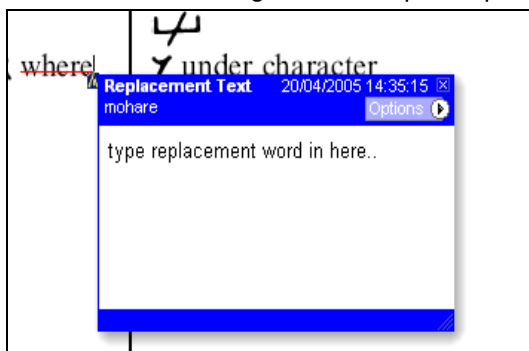


How to use it:

1. Right click into area of either inserted text or relevance to note
2. Select Add Note and a yellow speech bubble symbol and text box will appear
3. Type comment into the text box
4. Click the X in the top right hand corner of the note box to close.

Replacement text tool — For deleting one word/section of text and replacing it

Strikes red line through text and opens up a replacement text box.

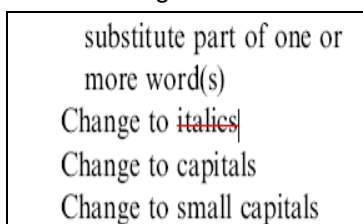


How to use it:

1. Select cursor from toolbar
2. Highlight word or sentence
3. Right click
4. Select Replace Text (Comment) option
5. Type replacement text in blue box
6. Click outside of the blue box to close

Cross out text tool — For deleting text when there is nothing to replace selection

Strikes through text in a red line.



How to use it:

1. Select cursor from toolbar
2. Highlight word or sentence
3. Right click
4. Select Cross Out Text

Approved tool — For approving a proof and that no corrections at all are required.

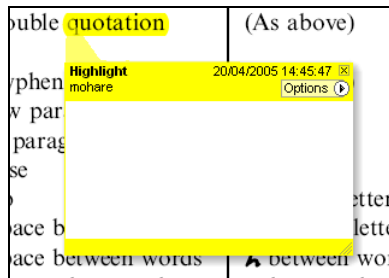


How to use it:

1. Click on the Stamp Tool in the toolbar
2. Select the Approved rubber stamp from the 'standard business' selection
3. Click on the text where you want to rubber stamp to appear (usually first page)

Highlight tool — For highlighting selection that should be changed to bold or italic.

Highlights text in yellow and opens up a text box.

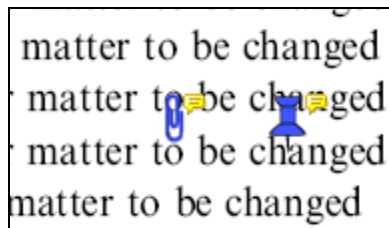


How to use it:

1. Select Highlighter Tool from the commenting toolbar
2. Highlight the desired text
3. Add a note detailing the required change

Attach File Tool — For inserting large amounts of text or replacement figures as a files.

Inserts symbol and speech bubble where a file has been inserted.

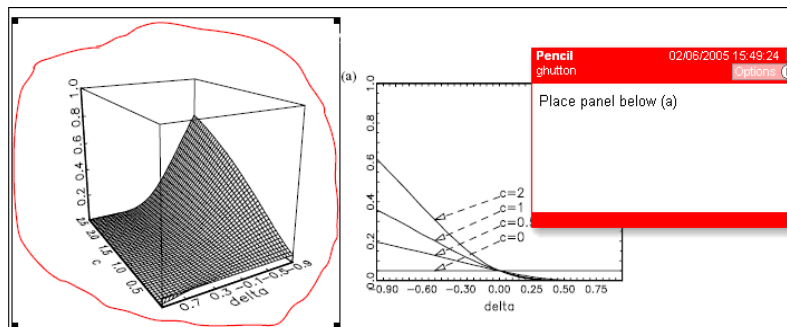


How to use it:

1. Click on paperclip icon in the commenting toolbar
2. Click where you want to insert the attachment
3. Select the saved file from your PC/network
4. Select appearance of icon (paperclip, graph, attachment or tag) and close

Pencil tool — For circling parts of figures or making freeform marks

Creates freeform shapes with a pencil tool. Particularly with graphics within the proof it may be useful to use the Drawing Markups toolbar. These tools allow you to draw circles, lines and comment on these marks.



How to use it:

1. Select Tools > Drawing Markups > Pencil Tool
2. Draw with the cursor
3. Multiple pieces of pencil annotation can be grouped together
4. Once finished, move the cursor over the shape until an arrowhead appears and right click
5. Select Open Pop-Up Note and type in a details of required change
6. Click the X in the top right hand corner of the note box to close.

Help

For further information on how to annotate proofs click on the Help button to activate a list of instructions:

