Learning synchronous context-free grammars with multiple specialised non-terminals for hierarchical phrase-based translation*

Felipe Sánchez-Martínez, Juan Antonio Pérez-Ortiz, Rafael C. Carrasco

Dep. de Llenguatges i Sistemes Informàtics, Universitat d'Alacant, Spain {fsanchez, japerez, carrasco}@dlsi.ua.es

Abstract

Translation models based on hierarchical phrase-based statistical machine translation (HSMT) have shown better performances than the non-hierarchical phrasebased counterparts for some language pairs. The standard approach to HSMT learns and apply a synchronous contextfree grammar with a single non-terminal. The hypothesis behind the grammar refinement algorithm presented in this work is that this single non-terminal is overloaded, and insufficiently discriminative, and therefore, an adequate split of it into more specialised symbols could lead to improved models. This paper presents a method to learn synchronous context-free grammars with a huge number of initial non-terminals, which are then grouped via a clustering algorithm. Our experiments show that the resulting smaller set of nonterminals correctly capture the contextual information that makes it possible to statistically significantly improve the BLEU score of the standard HSMT approach.

1 Introduction

Phrase-based statistical machine translation (PB-SMT) (Williams et al., 2016) has proven to be an effective approach to the task of machine translation. Even though in the last years neural systems have gained most of the attention from industry and academia, a number of recent works show that statistical approaches may still provide relevant results in hybrid, unsupervised or low-resource scenarios (Artetxe et al., 2018; Stahlberg et al., 2016). In PBSMT systems, the source-language (SL) sentence is split into non overlapping word sequences (known as *phrases*) and a translation into target-language (TL) is chosen for each phrase from a phrase table of bilingual phrase pairs extracted from a parallel corpus. A decoder efficiently chooses the splitting points and the corresponding TL equivalents by using the information provided by a set of features, which usually include probabilities provided by translation models as well as a language model. In spite of their performance, PBSMT systems are affected by some limitations due to the local strategy they follow. More specifically, they tend to overpass long range dependencies, which may negatively affect translation quality. Additionally, phrase reordering is usually instrumented by means of distortion heuristics and lexicalised reordering models or an operation-sequence model (Durrani et al., 2015) that cannot cope with the multiple structural divergences that are often necessary to translate between most language pairs.

Tree-based models address these issues by relying on a recursive representation of sentences that allows for *gaps* between words. Among these models, hierarchical phrase-based statistical machine translation (HSMT) (Chiang, 2005; Chiang, 2007) has gained lot of attention due to its relative simplicity and the lack of need for linguistic knowledge. HSMT infers synchronous contextfree grammars (SCFG); remarkably, HSMT infers grammars with a single non-terminal X.¹ Chiang (2007) also sets some additional restrictions on the extracted rules in order to contain the

^{*}Technical report, Transducens research group, Universitat d'Alacant.

¹Strictly speaking, two non-terminals are used since an additional non-terminal (used in the glue rules) is set as the initial symbol of the grammar (Chiang, 2005; Chiang, 2007).

combinatorial explosion in the number of rules, thus reducing decoding complexity.

This paper shows that the main limitation introduced in standard HSMT, namely, the use of a single non-terminal, can be overcome to improve translation quality. Our strategy differs from standard HSMT in that a different non-terminal is used for every possible gap when extracting the initial set of rules; this may easily result in millions of non-terminals which are then merged following an equivalence criterion, thus reducing the initial number of non-terminals in several orders of magnitude. Our experiments show that the resulting set of non-terminals correctly capture the contextual information that makes it possible to statistically significantly improve the BLEU score of the standard HSMT approach. However, additional research has to be carried out so that our method scales up with large training corpora.

Section 2 discusses related works. Section 3 then introduces the formalism of synchronous context-free grammars and the standard procedure to obtain them from parallel corpora and use them in HSMT. Section 4 presents our proposal for rule inference, including our criterion for variable equivalence and refinement. The experimental setup and the results of our experiments are presented in Section 5. Finally, the paper ends with some discussion and conclusions.

2 Related work

Probabilistic context-free grammars (PCFGs) have been traditionally used to model the generation of monolingual languages (Charniak, 1994). The refinement of probabilistic context-free grammars has been addressed before in a number of papers. For example, Matsuzaki et al. (2005) use a set of N latent symbols to annotate —and therefore specialize— non-terminals in a probabilistic context-free grammar. Every non-terminal A in the grammar is split into N non-terminals and the probabilities of the new productions are estimated to maximize the likelihood of the training set of strings.

In the approach for the refinement of PCFGs followed by Matsuzaki et al. (2005), the number of latent non-terminals must be small (N ranging from 1 to 16 in the experiments) since the training time and memory grow very fast with N. The input trees are also binarized before the estimation of the parameters to avoid an ex-

ponential growth in the number of productions contributing to every sum of probabilities. The method outperforms non-lexicalized approaches in terms of parsing accuracy measured over the Wall Street Journal (WSJ) portion of the Penn treebank (Marcus et al., 1993; Marcus et al., 1994).

In contrast, Petrov et al. (2006) apply a more sophisticated procedure for the specialization of the non-terminals in a PCFG. Instead of generating all possible annotations, their hierarchical splitting starts with the nearly one hundred tags in the WSJ corpus (binarization of the parse trees is also applied here) and leads to about one thousand symbols with a significant reduction in parsing error rate. The procedure splits iteratively every symbol into two sub-symbols and performs expectationmaximization optimization to estimate the probability of the productions. This method allows for a deeper recursive partition of some symbols (up to 6 times in the experiments with the WSJ corpus, as reported by the authors). Their results showed that significant improvements in the accuracy of parsing can be achieved with grammars which are more compact than those obtained by Matsuzaki et al. (2005). Parsing with the enhanced grammar can be accelerated with a coarse-to-fine scheme (Petrov and Klein, 2007) which prunes the items analyzed in the chart according to the estimations given by a simpler grammar where non-terminals are clustered into a smaller number of classes.

While strings are usually not enough to identify a particular PCFG, there exist methods which guaranteee convergence to the true grammar when structural information is available. For example, Pereira and Schabes (1992) apply expectation-maximization optimization to the identification of PCFGs from partially bracketed samples (sentences with parenthesis marking the constituent boundaries which the analysis should respect). The probabilistic tree grammars obtained after the optimization are non-deterministic (in the sense that multiple states are reachable as the result of processing a subtree bottom-up) and their trees must be binarized, a procedure that needs some linguistic guidance in order to generate meaningful probabilistic models.

The identification of probabilistic grammars has been addressed also through statemerging techniques. For example, such techniques have been applied to the case of regular grammars modelling string languages by Stolcke and Omohundro (1992) and Carrasco and Oncina (1999). In these methods, the initial grammar has one state or non-terminal per string-prefix in the sample and, then, the procedure looks for the optimal partition of the set of non-terminals, an approach that is similar to the one we also follow in this work. Convergence is guaranteed by keeping a frontier of states with all strings whose prefixes have been already examined.

Experimental work on the identification of regular string languages on sparse data (Lang et al., 1998) has shown the importance of exploring earlier the nodes about which there is more information available (a procedure known as *evidence-driven* merging). This result suggests that instead of using the straightforward breadth-first search (the simplest which guarantees convergence) one should consider more elaborate orders where those pairs of subtrees in the frontier with a higher number of observations are considered earlier.

Generalizing the merging techniques for string languages, Carrasco et al. (2001) define a procedure which identifies any regular tree language by comparing pairs of non-terminals: each initial nonterminal corresponds to a subtree in the sample and they are compared in a breadth-first mode and merged when they are found to be compatible.

As already commented in the introduction, Chiang (2005) extended the phrase-based approach in statistical machine translation to tree-based translation by adapting the rule algorithm to learn phrase pairs extraction with gaps (Koehn, 2010), and also presented a decoding method based on chart parsing, thus obtaining a working hierarchical phrasebased statistical machine translation system. Chiang (2007) then refined the method by introducing cube pruning to improve the efficiency of the chart decoder. Other authors (Maillette de Buy Wenniger and Sima'an, 2015; Vilar et al., 2010; Mylonakis and Sima'an, 2011) have introduced additional improvements to the rule learning procedure in the original proposal.

3 Synchronous context-free grammars in machine translation

Transductiongram-mars(Lewis II and Stearns, 1968)havebeen

used in SMT to model the hierarchical nature of translation which is implicit in the alignments between words in a bitext (Wu, 1997; Chiang, 2007). A bitext or bilingual text $B = (s_1, t_1)(s_2, t_2) \cdots (s_{|B|}, t_{|B|})$ consist of a finite sequence of sentence pairs, where every second component t_i , the target sentence, is the translation of the first component s_i in the pair, the source sentence.

A word-level alignment A(s,t) annotates every sentence pair (s,t) with a subset of $\{1,\ldots,|s|\} \times \{1,\ldots,|t|\}$ —where |s| and |t| are the sentence lengths— and provides those pairs of word positions in the source and target sentence which can be linked together according to a translation model.

grammar А transduction G(V, S, T, I, R, H) consists of a finite set of non-terminals $V = \{X_1, \ldots, X_N\}$, two sets of terminals —here, S and T consist of segments of words contained in source and target sentences respectively—, an initial symbol $I \in V$, a finite set of production rules (rules or productions, for short) $R = \{r_1, \dots, r_M\} \subset V \times (\mathcal{S} \cup V)^+ \times (\mathcal{T} \cup V)^+$ and a set of M one-to-one mappings $H = \{h_1, \dots, h_M\}$ with $h_m : \mathbb{N} \to \mathbb{N}$. For every rule $r_m = (X_n, \alpha, \beta) \in R$, h_m couples every instance of a non-terminal in α and an instance of the same non-terminal in β —therefore, α and β must have an identical number and type of non-terminals. A production $r_m = (X_n, \alpha, \beta)$ in a transduction grammar G will be written in the following as $X_n \rightarrow (\alpha, \beta)$ and their left and right components as $X_n = left(r_m)$ and $(\alpha, \beta) = \operatorname{right}(r_m)$, respectively.

The *synchronous context-free grammars* introduced by Chiang (2007) are transduction grammars whose productions have some restrictions:

- there is a single productive non-terminal X, in addition to the initial non-terminal I, that is, V = {I, X};
- 2. there is an upper limit on the length of the bilingual phrases of 10 words in either side;
- 3. there is an upper limit of 2 in the number of non-terminals that can appear in a production;
- 4. the size of the source side of production rules can have at most 5 terminals and non-terminals;

- 5. the productions cannot contain contiguous non-terminals on the source side and they must include some lexical content (terminals); and
- 6. two *glue rules* are defined to start derivations: $I \rightarrow (X, X)$ and $I \rightarrow (IX, IX)$.

The transduction grammars employed by Chiang (2007) restrict terminals to be in the set $\Phi(A, B)$ of bilingual phrase pairs obtained with the same extraction algorithm used in phrasebased SMT (Koehn, 2010, sect. 5.2.2): a *bilingual phrase pair* or *biphrase* is a pair in $S \times T$ which is consistent with the word alignments A provided by a statistical aligner for the bitext B.² The procedure also applies the following extraction rules:

- For every phrase pair (u, v) ∈ Φ(A, B), add a production X → (u, v) to R.
- For every production $r_m \in R$ such that $r_m = X \rightarrow (\alpha_1 u \alpha_2, \beta_1 v \beta_2)$ with $\alpha_i \in (S + V)^*$, $\beta_i \in (T+V)^*$, and $(u, v) \in \Phi(A, B)$, add the production $r_{M+1} = X \rightarrow (\alpha_1 X \alpha_2, \beta_1 X \beta_2)$ to R—with h_{M+1} extending h_m with a new link between the inserted X pair.

The previous procedure ends up generating production rules such as the following English–Chinese rule:

$$X \to (\text{hyu} X^{[1]} \text{you} X^{[2]} \mid \text{have} X^{[2]} \text{with} X^{[1]})$$

where the numbers in the superindexes are not used to represent different non-terminals but the coupling between non-terminals resulting from the corresponding one-to-one mapping h_i .

Each rule is given a probabilistic score. In order to find the most probable translation of an input sentence according to the grammar model, chart parsing is used at decoding time (Chiang, 2007).

4 A new method for grammar induction

Our strategy differs from the one by Chiang (2007) already introduced in the previous section in that a different non-terminal is used for every possible gap when extracting the initial set of rules; this may easily result in millions of non-terminals

which are then merged following an equivalence criterion, thus reducing the initial number of non-terminals in several order of magnitudes. Consequently, the following sections present the algorithm for extraction of production rules (Section 4.1), the criterion for considering two non-terminals as equivalent (Section 4.2) and the merging methods which join non-terminals based on the equivalence criterion (Section 4.3).

4.1 Extraction of production rules

The extraction phase assigns a different nonterminal to every production as follows (compare with the strategy proposed by Chiang (2007) and described in Section 3):

- 1. Start with the set of initial non-terminals $V = \{I\}$ I being the initial symbol— and empty set of production rules $R = \emptyset$.
- For every phrase pair (u, v) ∈ Φ(A, B), add a new non-terminal X_n to V —n being the current size of V—, and the new production X_n → (u, v), to R. If (u, v) is a sentence pair, then add also I → X_n to R. In contrast to Chiang (2007), the length of the phrasepairs used is not constrained. This results in a large number of production rules as well as in a large number of non-terminals but it is necessary to be able to reproduce each sentence pair in the training corpus.
- For every production r_k ∈ R such that r_k = X_i → (α₀α₁α₂, β₀β₁β₂) and there is a production X_j → (α₁, β₁) ∈ R, add the production r_{m+1} = X_i → (α₀X_jα₂, β₀X_jβ₂) to R —m being the size of R. Note that the subscript in the left-hand side is not changed, that is, left(r_{m+1}) = left(r_k). Note also that for every phrase pair (u, v) ∈ Φ(A, B) there is only one non-terminal X_n that can be derived to obtain (u, v), either by means of the immediate rule X_n → (u, v) or through a number of derivations starting with some other rule r_k ∈ R having left(r_k) = X_n.
- 4. Finally, the count of every non-terminal X_i is computed as done by Chiang (2007), that is, $C(X_i)$ is the number of occurrences in the training corpus of the phrase pair $(u, v) \in \Phi(A, B)$ generated by X_i . In order to generate the count for the production rule $c(r_k)$, the count $C(X_i)$ is equally distributed among

²Essentially, a pair (u, v) is a bilingual phrase pair if: u and v are segments in the source and target sentence, respectively; no word in u is aligned to a word not in v and vice versa; and at least one word in u is aligned to a word in v.

| Production | Count |
|---|---------------|
| $I \to (X_1, X_1)$ | 1 |
| $I \to (X_7, X_7)$ | 1 |
| $X_1 \rightarrow$ (das neue Haus, the new house) | $\frac{1}{6}$ |
| $X_2 \rightarrow$ (das neue, the new) | $\frac{1}{3}$ |
| $X_3 \rightarrow$ (neue Haus, new house) | $\frac{1}{3}$ |
| $X_4 \rightarrow (\text{das, the})$ | 2 |
| $X_5 \rightarrow (\text{neue, new})$ | 1 |
| $X_6 \rightarrow$ (Haus, house) | 2 |
| $X_7 \rightarrow$ (das Haus, the house) | $\frac{1}{3}$ |
| $X_1 \rightarrow (X_2 \text{ Haus, } X_2 \text{ house})$ | $\frac{1}{6}$ |
| $X_1 \rightarrow (\text{das } X_3, \text{the } X_3)$ | $\frac{1}{6}$ |
| $X_1 \rightarrow (X_4 \text{ neue Haus, } X_4 \text{ new house})$ | $\frac{1}{6}$ |
| $X_1 \rightarrow$ (das X_5 Haus, the X_5 house) | $\frac{1}{6}$ |
| $X_1 \rightarrow$ (das neue X_6 , the new X_6) | $\frac{1}{6}$ |
| $X_2 \rightarrow (X_4 \text{ neue, } X_4 \text{ new})$ | $\frac{1}{3}$ |
| $X_2 \rightarrow (\text{das } X_5, \text{the } X_5)$ | $\frac{1}{3}$ |
| $X_3 \rightarrow (X_5 \text{ Haus, } X_5 \text{ house})$ | $\frac{1}{3}$ |
| $X_3 \rightarrow (\text{neue } X_6, \text{new } X_6)$ | $\frac{1}{3}$ |
| $X_7 \rightarrow (X_4 \text{ Haus}, X_4 \text{ house})$ | $\frac{1}{3}$ |
| $X_7 \rightarrow (\text{das } X_6, \text{the } X_6)$ | $\frac{1}{3}$ |

Table 1: Productions and fractional counts for the monotonic alignment of the sentence pairs ("das neue Haus", "the new house") and ("das Haus", "the house"). The first group of rules includes the *glue* rules; the second group includes those rules added in step 2 of the algorithm in Section 4.1; the third group is made up of the rules added in step 3. The fractional production counts, computed as described in step 4, are shown in the second column.

all the productions r_k generating that phrase pair, that is, between $X_i \rightarrow (u, v)$ and the other productions with X_i in the left-hand side added in step 3.

Figure 1 shows the resulting production rules and their fractional counts after applying our extraction procedure to the German–English sentence pairs ("das neue Haus", "the new house") and ("das House", "the house"), assuming a monotonic alignment in which the *i*-th word of one sentence is aligned with the *i*-th word of the other sentence. Following Chiang (2007), rules with no lexical content, such as $X_1 \rightarrow (X_2X_3, X_3X_2)$, have not been considered.

4.2 Determining equivalent non-terminals

As will be presented in Section 4.3, our method will merge pairs of equivalent non-terminals. We will denote with $X_i \sim X_j$ the fact that X_i and X_j are equivalent and, thus, they must be merged. Then, $X_i \sim X_j$ implies that for all pairs of productions r_m and r_n which, for some α and β in $(V \cup T)^+$ and X_k and X_l in V, have the form

$$r_m = X_k \to \alpha X_i \beta$$

and

$$r_n = X_l \to \alpha X_j \beta$$

the following equality is (approximately) satisfied

$$\frac{c(r_m)}{C(X_i)} \approx \frac{c(r_n)}{C(X_j)} \tag{1}$$

and, recursively, $X_k \sim X_l$. The approximate matching between the above quotients is probabilistic in nature and must be defined therefore in terms of a stochastic test, such as the Hoeffding (1963) bound. When using this test for proportion comparison, two proportions c_1/C_1 and c_2/C_2 are not statistically different if:

$$\left|\frac{c_1}{C_1} - \frac{c_2}{C_2}\right| < \sqrt{\frac{-\log\frac{\alpha}{2}}{2\frac{C_1C_2}{C_1 + C_2}}}$$

where α is the confidence level.³ Note that the possible use of α in the proportion test is twofold: on the one hand, we may set α to a fixed value in order to test whether two proportions are statistically different; on the other hand, we may compute the value of α for which the test changes from true to false and use this value as a continuous measure of non-terminal *dissimilarity*: after isolating α in the Hoeffding test equation and removing terms which are constant across different evaluations, we can easily arrive to a function *D* that provides the dissimilarity of two non-terminals in a particular comparison context based on their respective counts:

$$D(C_1, C_2, c_1, c_2) = \frac{C_1 C_2}{C_1 + C_2} \left(\frac{c_1}{C_1} - \frac{c_2}{C_2}\right)^2$$
(2)

The dissimilarity of two variables is therefore obtained as the *maximum* value of D for all the contexts representing the different tests performed upon the variables as described in the previous algorithm.

³Following Habrard et al. (2003), we use a Fisher exact test as a back-off test in the experiments when the number of observations is small.

4.2.1 Example of equivalence computation

In order to gain some insight on the meaning of equivalence between non-terminals, let us now consider, for instance, the comparison between X_3 and X_6 in the example in Table 1, which can be considered plausible candidates for equivalence since they both generate a noun phrase pair. The fractional number of occurrences for X_3 and X_6 are given by $C(X_3) = \frac{1}{3} + \frac{1}{3} + \frac{1}{3} = 1$ and $C(X_6) = 2$. Note that $C(X_6)$ receives contributions from both sentence pairs.

Non-terminal X_3 is only in $X_1 \rightarrow (\text{das } X_3, \text{ the } X_3)$ —with weight $\frac{1}{6}$ — and therefore, equivalence implies that a production with right-hand side (das X_6 , the X_6) must be in R with a weight which is consistent with Equation (1). Indeed, $X_7 \rightarrow (\text{das } X_6, \text{ the } X_6)$ is in R with weight $\frac{1}{3}$: if $X_3 \sim X_6$ both production will appear with a relative frequency which must be similar to the relative frequency for X_3 and X_6 . In this case,

$$\frac{c(X_7 \to (\text{das } X_6, \text{ the } X_6))}{C(X_6)} = \frac{1/3}{2} = \frac{1}{6}$$

while

$$\frac{c(X_1 \to (\text{das } X_3, \text{ the } X_3))}{C(X_3)} = \frac{1/6}{1} = \frac{1}{6}$$

Furthermore, the left-hand sides, X_1 and X_7 must be also equivalent.⁴

The quotients above reflect that the word pair (Haus, house) has been observed in two different contexts: after the biword (das, the) and also following (neue, new); in contrast, the phrase pair (neue Haus, new house) has been only observed after (das, the). This asymmetry will lead to different values in the relative frequencies. Of course, one cannot expect to draw definitive conclusions from such a tiny bitext —clearly, a much larger sample will be needed to extract reliable estimates but the example illustrates how the frequency estimates provide hints to differentiate between equivalent non-terminal pairs and those which, in contrast, should remain distinct.

Once productions with X_3 on the right-hand side have been checked, those with X_6 should be checked although, by the symmetry of the test, only those that have no correspondent X_3 -production. In our example, the tests will be

$$\frac{c(X_1 \to (\text{das neue } X_6, \text{ the new } X_6)}{C(X_6)} \approx 0$$

and

$$\frac{c(X_3 \to (\text{neue } X_6, \text{ new } X_6)}{C(X_6)} \approx 0$$

and, thus, $X_1 \sim X_3$.

Note that, if done carefully, recursion is always finite because the comparison between nonterminals is consistent with the depth of the subtree associated to each non-terminal.

An efficient algorithm has to be carefully designed in order to avoid duplicate calls when the equivalence between X_i and X_j is tested (since equivalence is a symmetric and transitive relation).

4.3 Merging variables

Given the equivalence criterion presented in the previous section, an algorithm for non-terminal merging can be run in order to group equivalent non-terminals and reduce the initial number. We have evaluated two different algorithms: an adaptation of the Blue-Fringe algorithm and a k-medoids clustering algorithm.

4.3.1 The Blue-Fringe algorithm

The following description, based on that by Lang et al. (1998) of the procedure proposed by Juillé and Pollack (1998), adapts the Blue-Fringe algorithm for the identification of regular string languages to the case of transduction grammars.

The procedure splits the set of non-terminals V into three subsets: red (the kernel K of mutually non-equivalent non-terminals), blue (the frontier F being explored) and white (the subset W = V - K - F of non-terminals with pending classification). At every iteration a non-terminal is removed from the frontier F and either it becomes a new member of K or it is merged with an equivalent one in K (and, thus, removed from V). As will be seen immediately, after every addition or merge, some non-terminals in W can be moved to F.

In order to avoid a possible infinite recursion in equivalence tests, non-terminals in F must not produce any non-terminal in K through derivation. This condition can be guaranteed if a nonterminal X_n can only enter F if all the content on the right-hand side of productions with the form $X_n \to (\alpha, \beta)$ is either lexical or already in K: this implies that any production r_m such that $X_n \in F$ is in right(r_m) satisfies left(r_m) $\in W$.

⁴Since they only appear in *I*-productions, this is trivially true.

Initially, leaves (non-terminals producing only irreducible phrase pairs) are red; non-terminals with a production with only lexical content and red non-terminals are blue; all other non-terminals are white. Our method will merge pairs of equivalent non-terminals by comparing those with a higher number of observations first (Lang et al., 1998). The policy described by Juillé and Pollack (1998) performs the following actions while there are still blue non-terminals:

- 1. Evaluate all red-blue merges.
- 2. If there exists a blue non-terminal that cannot be merged with any red non-terminal (meaning that no equivalent non-terminal is found), promote one of the shallowest such blue nonterminals to red (ties are broken at random).
- 3. Otherwise (if no blue non-terminal can be promoted), perform the red-blue merge with highest score. This score here is be based on the fractional number of occurrences of the corresponding non-terminals.

Then, some white non-terminals are moved to the frontier F as stated before.

4.3.2 *k*-medoids clustering

The k-medoids algorithm is a clustering algorithm similar to the well-known k-means, but with the particularity that the center of each cluster (the *medoid*) is a point in the data, which is specially relevant in our case since the representative of each cluster must be an existing non-terminal in the grammar. Unlike the Blue-Fringe algorithm, which attains a different number of final non-terminals depending on the confidence level α , the parameter set *a priori* in this case is the number of final clusters (i.e. non-terminals) k. Note that when following this merge approach, the dissimilarity between non-terminals in Equation (2) is used to compute the required distances.

5 Experimental setup

We trained and evaluated our grammar induction procedure on the English–Spanish language pair using a small fraction of the EMEA corpus.⁵ Table 2 provide additional information about the corpora used in the experiments.

| Corpus | # Sentences | # Words (en/es) |
|--------|-------------|-------------------|
| train | 73,372 | 519,763 / 556,453 |
| dev | 2,000 | 22,410 / 25,219 |
| test | 3,000 | 33,281 / 37,492 |

Table 2: Number of sentences and words in each language for the corpora used in the experiments.

In order to have a manageable initial set of non-terminals and productions and make the problem computationally affordable we limited the sentences to be included in the training corpus to a maximum of 20 words. In addition, instead of using the words themselves when defining the initial set of non-terminals we used word classes. In particular, we used 10 word classes obtained by running mkcls⁶ (Och, 1999) for 5 iterations. As a result, the initial set of non-terminals contains 852,423 non-terminals and the amount of productions, not including those involving the initial nonterminal I, is 2,055,902.

All the experiments were carried out with the free/open-source SMT system (Koehn et al., 2007), Moses release 2.1.1. GIZA++ (Och and Ney, 2003) was used for computing words alignments. KenLM was used to train a 5-gram language model on a monolingual corpus made of Europarl v7,⁷ News Commentary $v8^8$ and the EMEA sentences in the training corpus; in total the corpus used for training the language model consists of 3,423,702 Spanish sentences. The weights of the different feature functions were optimised by means of minimum error rate training (Och, 2003). The parallel corpora were tokenised and truecased before training, as were the development and test sets used.

6 Results

Table 3 reports the BLEU scores obtained on the test set when running the Blue-Fringe algorithm for different values of α . The amount of output non-terminals in the inferred context-free grammar is also reported. The best results are obtained with $\alpha = 10^{-2}$. The performance of the baseline hierarchical phrase-based system (Chiang, 2005) is 0.5818. The difference in performance is statisti-

⁵http://opus.nlpl.eu/EMEA.php

⁶http://www.statmt.org/moses/giza/mkcls.html
⁷http://www.statmt.org/wmt13/

training-monolingual-europarl-v7.tgz

⁸http://www.statmt.org/wmt13/

training-monolingual-nc-v8.tgz

| α | # final non-terminals | BLEU |
|-----------|-----------------------|--------|
| 10^{-1} | 4,434 | 0.5838 |
| 10^{-2} | 2,346 | 0.5868 |
| 10^{-3} | 1,606 | 0.5859 |
| 10^{-4} | 1,261 | 0.5855 |
| 10^{-5} | 1,074 | 0.5866 |

Table 3: BLEU scores obtained by the Blue-Fringe clustering algorithm for different values of α . The performance of the baseline HSMT system is 0.5818.

| # clustered | # final non-terminals k | | |
|---------------|---------------------------|--------|--------|
| non-terminals | 2 | 3 | 4 |
| 125 | 0.5975 | 0.5991 | 0.5952 |
| 250 | 0.5986 | 0.6000 | 0.5942 |
| 500 | 0.5946 | 0.5972 | 0.5936 |
| 1000 | 0.5950 | 0.5998 | 0.5939 |
| 2500 | 0.5985 | 0.5984 | 0.5964 |
| 5000 | 0.5947 | 0.5991 | 0.5947 |

Table 4: BLEU scores obtained by the k-medoids clustering method for different sizes of the subset of non-terminals over which the clustering is performed and for different numbers of clusters (i.e. final non-terminals). The performance of the baseline HSMT system is 0.5818.

cally significant for the figures in bold according to paired bootstrap resampling (Koehn, 2004) with p = 0.05.

Table 4 shows the BLEU scores obtained on the test set when the k-medoids clustering algorithm is run over the n most frequent non-terminals (# clustered non-terminals) to obtained a pre-fixed number of clusters, that is, of non-terminals in the inferred grammar; the remaining non-terminals are then added to the nearest cluster after the algorithm finishes. The table reports results for 2, 3 and 4 clusters; although we tried with more clusters these are the numbers of clusters for which we got the best results. The difference in performance with the baseline is statistically significant for all the figures reported in the table according to paired bootstrap resampling (Koehn, 2004) with p = 0.05.

The results obtained when the number of nonterminals over which the k-medoids is run is set to 250 and the number of cluster to obtained is set to 3 are better than those obtained with the Blue-Fringe algorithm and better that the results achieved by the baseline. The k-medoids allow us to get an 3% improvement in BLEU with just three non-terminals, in contrast with the thousand nonterminals obtained by the Blue-Fringe algorithm.

7 Conclusions

This work extends the well-known algorithm for rule extraction in hierarchical statistical machine translation originally proposed by Chiang (2007). Our proposal allows for more than one nonterminal in the resulting synchronous context-free grammar, thus incorporating a specialisation in the resulting non-terminals. Our method works by initially creating a different non-terminal for every possible gap when extracting the initial set of rules; this may easily result in millions of non-terminals which are then merged following a novel equivalence criterion for non-terminals. Two merging strategies are proposed: one inspired on the Blue-Fringe algorithm that joins non-terminals on a one-by-one basis, and another one that performs a k-medoids clustering over a reduce set of nonterminals. A statistically significant improvement in BLEU as with respect to the original method is obtained with both merging criteria.

For the experiments we used a small parallel corpus and restricted the length of the parallel sentences in the training corpus to 20 words. This was necessary in order to be able to run the merging algorithm in reasonable time. Recall that, in contrast to (2007), we do not restrict the length of the phrase pairs in order to be able to reproduce the parallel sentences in the training corpus; otherwise long-range reorderings happening near the root of the parse tree of the sentences would not be possible. We tried different methods for filtering the rule table before applying the approach described in this paper so as to be able to use larger corpora and longer sentences, although with no success. A deeper exploration of potential optimizations is necessary.

Acknowledgements Work supported by the Spanish government through project EFFOR-TUNE (TIN2015-69632-R). The authors thank Mikel L. Forcada for his helpful comments.

References

- [Artetxe et al.2018] Artetxe, Mikel, Gorka Labaka, and Eneko Agirre. 2018. Unsupervised statistical machine translation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*.
- [Carrasco and Oncina1999] Carrasco, Rafael C. and José Oncina. 1999. Learning deterministic regular grammars from stochastic samples in polyno-

33(1):1-20.

- [Carrasco et al.2001] Carrasco, Rafael C., José Oncina, and Jorge Calera-Rubio. 2001. Stochastic inference of regular tree languages. Machine Learning, 44(1/2):185-197.
- [Charniak1994] Charniak, Eugene. 1994. Statistical Language Learning. MIT Press.
- [Chiang2005] Chiang, David. 2005. A hierarchical phrase-based model for statistical machine translation. In Proc. ACL, pages 263-270.
- [Chiang2007] Chiang, David. 2007. Hierarchical phrase-based translation. Computational Linguistics, 33(2):201-228.
- [Durrani et al.2015] Durrani, Nadir, Helmut Schmid, Alexander Fraser, Philipp Koehn, and Hinrich Schutze. 2015. The operation sequence model combining n-gram-based and phrase-based statistical machine translation. Computational Linguistics, 41(2):157-186.
- [Habrard et al.2003] Habrard, Amaury, Marc Bernard, and Marc Sebban. 2003. Improvement of the state merging rule on noisy data in probabilistic grammatical inference. In Proceedings of ECML 2003, 14th European Conference on Machine Learning, pages 169-180.
- [Hoeffding1963] Hoeffding, Wassily. 1963. Probability inequalities for sums of bounded random variables. Journal of the American Statistical Association, 58(301):13-30.
- [Juillé and Pollack1998] Juillé, Hugues and Jordan B. Pollack. 1998. A sampling-based heuristic for tree search applied to grammar induction. In Mostow, Jack and Chuck Rich, editors, Proceedings of the Fifteenth National Conference on Artificial Intelligence and Tenth Innovative Applications of Artificial Intelligence Conference, AAAI 98, IAAI 98, pages 776-783. AAAI Press / The MIT Press.
- [Koehn et al.2007] Koehn, Philipp, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondřej Bojar, Alexandra Constantin, and Evan Herbst. 2007. Moses: Open source toolkit for statistical machine translation. In Proceedings of the 45th Annual Meeting of the ACL on Interactive Poster and Demonstration Sessions, ACL '07, pages 177-180. Association for Computational Linguistics.
- [Koehn2004] Koehn, P. 2004. Statistical significance tests for machine translation evaluation. In Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing, pages 388-395.
- [Koehn2010] Koehn, Philipp. 2010. Statistical machine translation. Cambridge University Press.

- mial time. Informatique Théorique et Applications, [Lang et al. 1998] Lang, Kevin J., Barak A. Pearlmutter, and Rodney A. Price. 1998. Results of the Abbadingo One DFA learning competition and a new evidence-driven state merging algorithm. In Proceedings of the 4th International Colloquium on Grammatical Inference, ICGI '98, pages 1-12. Springer-Verlag.
 - [Lewis II and Stearns1968] Lewis II, P. M. and R. E. Stearns. 1968. Syntax-directed transduction. Journal of the ACM, 15(3):465-488.
 - [Maillette de Buy Wenniger and Sima'an2015] Maillette de Buy Wenniger, Gideon and Khalil Sima'an. 2015. Labeling hierarchical phrase-based models without linguistic resources. Machine Translation, 29(3):225-265.
 - [Marcus et al. 1993] Marcus, Mitchell P., Mary Ann Marcinkiewicz, and Beatrice Santorini. 1993. Building a large annotated corpus of english: the Penn Treebank. Comput. Linguist., 19(2):313-330, June.
 - [Marcus et al.1994] Marcus, Mitchell P., Grace Kim, Mary Ann Marcinkiewicz, Robert MacIntyre, Ann Bies, Mark Ferguson, Karen Katz, and Britta Schasberger. 1994. The Penn Treebank: Annotating predicate argument structure. In Human Language Technology, Proceedings of a Workshop held at Plainsboro, New Jerey, USA, March 8-11, 1994. Morgan Kaufmann.
 - [Matsuzaki et al.2005] Matsuzaki, Takuya, Yusuke Miyao, and Jun'ichi Tsujii. 2005. Probabilistic cfg with latent annotations. In Knight, Kevin, Hwee Tou Ng, and Kemal Oflazer, editors, ACL 2005, 43rd Annual Meeting of the Association for Computational Linguistics, Proceedings of the Conference, 25-30 June 2005, University of Michigan, USA. The Association for Computer Linguistics.
 - [Mylonakis and Sima'an2011] Mylonakis, Markos and Khalil Sima'an. 2011. Learning hierarchical translation structure with linguistic annotations. In Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies, pages 642-652.
 - [Och and Ney2003] Och, Franz Josef and Hermann Ney. 2003. A systematic comparison of various statistical alignment models. Comput. Linguist., 29(1):19-51.
 - [Och1999] Och, Franz Josef. 1999. An efficient method for determining bilingual word classes. In Ninth Conf. of the Europ. Chapter of the Association for Computational Linguistics, EACL'99, pages 71-76, June.
 - [Och2003] Och, Franz Josef. 2003. Minimum error rate training in statistical machine translation. In Proceedings of the 41st Annual Meeting on Association for Computational Linguistics, pages 160-167.

- [Pereira and Schabes1992] Pereira, Fernando C. N. and Yves Schabes. 1992. Inside-outside reestimation from partially bracketed corpora. In Thompson, Henry S., editor, *Proceedings 30th Annual Meeting of the Association for Computational Linguistics*, pages 128–135.
- [Petrov and Klein2007] Petrov, Slav and Dan Klein. 2007. Improved inference for unlexicalized parsing. In Sidner, Candace L., Tanja Schultz, Matthew Stone, and ChengXiang Zhai, editors, *Proceedings Human Language Technology Conference of the North American Chapter of the Association of Computational Linguistics*, pages 404–411. The Association for Computational Linguistics.
- [Petrov et al.2006] Petrov, Slav, Leon Barrett, Romain Thibaux, and Dan Klein. 2006. Learning accurate, compact, and interpretable tree annotation. In Calzolari, Nicoletta, Claire Cardie, and Pierre Isabelle, editors, *Proceedings ACL 2006, 21st International Conference on Computational Linguistics and* 44th Annual Meeting of the Association for Computational Linguistics. The Association for Computer Linguistics.
- [Stahlberg et al.2016] Stahlberg, F., E. Hasler, A Waite, and B Byrne. 2016. Syntactically guided neural machine translation. In Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, pages 299–305.
- [Stolcke and Omohundro1992] Stolcke, Andreas and Stephen M. Omohundro. 1992. Hidden Markov Model induction by bayesian model merging. In Hanson, Stephen Jose, Jack D. Cowan, and C. Lee Giles, editors, Advances in Neural Information Processing Systems 5, pages 11–18. Morgan Kaufmann.
- [Vilar et al.2010] Vilar, David, Daniel Stein, Stephan Peitz, and Hermann Ney. 2010. If I only had a parser: Poor man's syntax for hierarchical machine translation. In *Proc. of the Int. Workshop on Spoken Language Translation (IWSLT)*.
- [Williams et al.2016] Williams, Philip, Rico Sennrich, Matt Post, and Philipp Koehn. 2016. *Syntax-based statistical machine translation*. Morgan & Claypool Publishers.
- [Wu1997] Wu, Dekai. 1997. Stochastic inversion transduction grammars and bilingual parsing of parallel corpora. *Computational Linguistics*, 23(3):377–403.