

Grammar based statistical MT on Hadoop
An end-to-end toolkit for large scale PSCFG based MT

Ashish Venugopal, Andreas Zollmann

Abstract

This paper describes the open-source Syntax Augmented Machine Translation (SAMT)¹ on Hadoop toolkit—an end-to-end grammar based machine statistical machine translation framework running on the Hadoop implementation of the MapReduce programming model. We present the underlying methodology of the SAMT approach with detailed instructions that describe how to use the toolkit to build grammar based systems for large scale translation tasks.

1. Introduction

1.1. PSCFG approaches to Machine Translation

Syntax Augmented Machine Translation (SAMT) (Zollmann and Venugopal, 2006) defines a specific parameterization of the probabilistic synchronous context-free grammar (PSCFG) approach to machine translation. PSCFG approaches take advantage of nonterminal symbols, as in monolingual parsing, to generalize beyond purely lexical translation. Consider the example rule below:

$$@VP \rightarrow ne @VB_1 pas \# do not @VB_1 : w$$

representing the discontinuous translation of the French words “ne” and “pas” to “do not”, in the context of the labeled nonterminal symbol “@VB” (representing the syntactic constituent type of Verb). These rules seem considerably more complex than weighted word-to-word

¹Released under the GNU Lesser General Public License, version 2

rules (Brown et al., 1993), or phrase-to-phrase rules (Koehn, Och, and Marcu, 2003, Och and Ney, 2004) but can be viewed as natural extensions to these well established approaches. An introduction to PSCFG approaches to machine translation can be found in (Chiang and Knight, 2006).

(Chiang, 2005) describes a procedure to learn PSCFG rules from word-aligned parallel corpora, using the phrase-pairs from (Koehn, Och, and Marcu, 2003) as a lexical basis for the grammar. SAMT (Zollmann and Venugopal, 2006) extends the procedure from (Chiang, 2005) to assign labels to nonterminal symbols based on target language phrase structure parse trees.

In this paper, we describe an end-to-end statistical machine translation framework—SAMT on Hadoop—to learn and estimate parameters for PSCFG grammars from word-aligned parallel corpora (*training*), and perform translation (*decoding*) with these grammars under a log-linear translation model (Och and Ney, 2004). While our framework specifically implements (Chiang, 2005) and (Zollmann and Venugopal, 2006), the training and decoding algorithms in our toolkit can be easily replaced to experiment with alternative PSCFG parameterizations like (Galley et al., 2006, Wu, 1997) or alternative decoding approaches such as (Petrov, Haghghi, and Klein, 2008, Zhang and Gildea, 2008). The algorithms in this toolkit are implemented upon Hadoop (Cutting and Baldeschwieler, 2007), an open-source implementation of the MapReduce (Dean and Ghemawat, 2004) framework, which supports distribution computation on large scale data using clusters of commodity hardware. We report empirical results that demonstrate the use of the SAMT toolkit on large scale translation tasks.

1.2. The SAMT toolkit

Our toolkit, when used in concert with other open-source components and publicly available corpora, contains all of the necessary components to build and evaluate grammar based statistical machine translations systems. The primary components of the toolkit are listed below:

- A top level push-button script that provides experimental work-flow management and submits jobs to the underlying Hadoop framework.
- Components to build and estimate parameters for the grammars described in (Chiang, 2005) and (Zollmann and Venugopal, 2006).
- Tools to filter large translation grammars and n-gram language models to build small sentence specific models that can be easily loaded into memory during decoding.
- A bottom-up dynamic chart parsing decoder based on (Chappelier and Rajman, 1998) which supports grammars with more than 2 nonterminals symbols per rule. The decoder outputs n-best lists with optional annotations that facilitate discriminative training.
- An implementation of Minimum Error Rate (MER) training (Och, 2003), extended to perform feature selection.

The SAMT toolkit requires the following inputs that are easily generated by existing open-source tools.

- Word aligned parallel corpora. For small resource tasks, word-alignments can be gen-

erated using the GIZA++ toolkit (Och and Ney, 2003), while large-resource tasks can be aligned using (Dyer et al., 2008), a parallelized GIZA++ implementation on MapReduce.

- (Zollmann and Venugopal, 2006) requires target language parse trees for each sentence in the training data. SAMT on Hadoop interfaces to the parser from (Charniak, 2000) to parse the target side of the parallel corpora on Hadoop.
- N-Gram language models built via the SRILM toolkit (Stolcke, 2002) are used as features during decoding.

1.3. SAMT on Hadoop

As discussed in (Dyer et al., 2008), there are significant computational challenges in estimating the component models of a statistical machine translation system from the large parallel and monolingual corpora that yield state-of-the-art translation quality. (Dyer et al., 2008) apply the MapReduce (Dean and Ghemawat, 2004) programming model to distribute the estimation of word-alignment models (Brown et al., 1993) on a cluster of commodity machines. The MapReduce model requires that large computational tasks be split into two distinct steps, a Map step and a Reduce step. In the Map step, input data is processed by parallel tasks generating intermediate output in the form of key-value pairs. Map tasks tend to be those aspects of the overall computation that can be performed with access to only a limited, arbitrary portion of the input data. In the Reduce step, tasks running in parallel take as input intermediate Map output, with the guarantee that a single Reduce task will receive all intermediate key-value pairs that share the same key.

The SAMT toolkit is built upon Hadoop (Cutting and Baldeschwieler, 2007), an open-source implementation of the MapReduce model to distribute the estimation of PSCFG grammars and to perform decoding. Training and decoding are broken up into a series of MapReduce tasks, called *phases*, which are performed sequentially, transforming input data into a PSCFG grammar, and using the grammar to translate development and test sentences. Phase outputs are stored on the Hadoop Distributed File System (HDFS), a highly fault tolerant file system that is accessible by all cluster machines. Most SAMT phases are run sequentially, using output from previous phases as input.

1.4. Running SAMT on Hadoop

Detailed instructions for downloading and building the SAMT package are available at the toolkit’s website². The top-level script that is used to build PSCFG grammars and perform translation in SAMT is `mr_runner.pl`; this script is responsible for interpreting user parameter files, submitting jobs to the Hadoop infrastructure and checking for error codes from Hadoop jobs. For each phase in the SAMT pipeline, this script generates a Hadoop-on-Demand script file, which is interpreted by the Hadoop architecture to submit jobs to a cluster

²www.cs.cmu.edu/~zollmann/samt

of dedicated machines³. Each script file refers to pre-compiled Map and Reduce binaries and specifies the input and output paths for the phase. With respect to a machine translation workflow, this script is ultimately responsible for building PSCFG grammars, running MER training and translating unseen test data, and allows output from previous experiments to initialize new experiments.

The SAMT toolkit is distributed with parameter files in the `examples` directory that can be used to re-generate published results from (Zollmann, Venugopal, and Vogel, 2008). In the remainder of this paper, we describe the SAMT methodology and important user parameters in our toolkit that impact translation quality and runtime. For a more formal description of the individual MapReduce phases in the SAMT pipeline, see (Zollmann, Venugopal, and Vogel, 2008).

2. Syntax Augmented Machine Translation

2.1. Phrase and SAMT Rule Extraction

In this section, we describe Syntax Augmented Machine Translation (SAMT) (Zollmann and Venugopal, 2006), a specific instantiation of the PSCFG formalism that is implemented in the SAMT on Hadoop toolkit. SAMT extends the purely hierarchical grammar proposed in (Chiang, 2005) to use nonterminal labels learned from target language parse trees. The inputs to the SAMT rule extraction procedure are tuples, $\langle f, e, Phrases(a, f, e), \pi \rangle$, where f is a source sentence, e is a target sentence, a is a word-to-word alignment associating words in f with words in e , $Phrases(a, e, f)$, are the set of phrase pairs (source and target phrases) consistent with the alignment a (Koehn, Och, and Marcu, 2003, Och and Ney, 2004), and π is a phrase structure parse tree of e . SAMT rule extraction associates each phrase pair from $Phrases(a, e, f)$ with a left-hand-side label, and then applies the rule extraction procedure from (Chiang, 2005) to generate rules with *labeled* nonterminal symbols.

Consider the example alignment graph (a word alignment and target language parse tree as defined in (Galley et al., 2006)) for the example French-to-English sentence in Figure 1. The phrase extraction method from (Koehn, Och, and Marcu, 2003), extracts all phrase pairs where no word inside the phrase pair is aligned to a word outside the phrase pair. The following phrase-pairs, (with source and target sides separated by the “#” symbol) would be extracted for our example sentence:

```

il # he
va # go
ne va pas # does not go
ne va pas # not go
il ne va pas # he does not go

```

³While these scripts assume the Hadoop-on-Demand machine requisitioning model, the toolkit can be easily modified to submit jobs to a single global machine pool

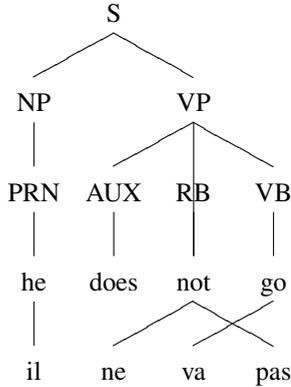


Figure 1: Alignment graph (word alignment and target parse tree) for a French-English sentence pair.

Phrase Extraction is the first phase of the SAMT toolkit, annotating each sentence-pair of the training corpus with a set of phrase pairs extracted from that sentence pair. We use a single toolkit binary: **MapExtractPhrases**, run as Hadoop Map step (there is no Reduce step in this phase). This binary takes a single numerical argument which determines the maximum length of the initial phrase extracted from word-aligned data. This limit has an impact on the size and nature of the final grammar. Typically, phrase limits are significantly smaller than the length of the parallel sentence, preventing very long distance reordering effects from being captured in the grammar.

The next phase includes rule extraction (Map step, binary **MapExtractRules**) on a per-sentence basis, and merging and counting of identical rules (Reduce step, binary **MergeRules**). SAMT assigns a left-hand-side (lhs) label to every phrase pair extracted from the current sentence-pair, based on the corresponding target language parse tree π , forming *initial rules*. These labels are assigned based on the constituent spanning the target side word sequence in π . When the target side of the phrase-pair is spanned by a single constituent in π , the constituent label is assigned as the lhs for the phrase pair. If the target side of the phrase is not spanned by a single constituent in π , we use the labels of subsuming, subsumed, and neighboring constituents in π to assign an extended label of the form $C_1 + C_2$, C_1/C_2 , or $C_2 \setminus C_1$ (similar in motivation to the labels in (Steedman, 1999)), indicating that the phrase pair’s target side spans two adjacent syntactic categories (e.g., *she went*: $NP+VB$), a partial syntactic category C_1 missing a C_2 at the right (e.g., *the great*: NP/NN), or a partial C_1 missing a C_2 at the left (e.g., *great wall*: $DT \setminus NP$), respectively. The label assignment is attempted in the order just described, i.e., assembling labels based on ‘+’ concatenation of two subsumed constituents is preferred, as smaller constituents tend to be more accurately labeled. If no label is assignable by either of these three methods, and the parameter ‘-allow_double_plus 1’ is set, we try triple-concatenation to create a label of the form $C_1 + C_2 + C_3$. If this approach

do not yield a label or if ‘-allow_double_plus 0’, a default label ‘_FAIL’ is assigned.

An ambiguity arises when unary rules $N_1 \rightarrow \dots \rightarrow N_m$ in the target parse tree are encountered, such as the NP→PRN subtree in Figure 1. Depending on the parameter ‘-unary_category_handling’, we use the bottom-most label (parameter value ‘bottom’), the top-most (‘top’), or a combined label $N_m : \dots : N_1$ (‘all’, this is the default).

An alternative method of assigning labels to phrase pairs can be activated by specifying the parameter ‘-use_only_pos’. In this variant, labeling is performed merely based on the part-of-speech (POS) tags of the first word $POS1$ and last word $POS2$ of the target phrase, resulting in the label ‘ $POS1-POS2$ ’. In general, the SAMT approach can take advantage of any labeling techniques that assigns labels to arbitrary initial phrase pairs. Alternative techniques could include using source language constituent labels, or automatically induced labels.

The following initial rules would be extracted for our example sentence pair with default parameter settings:

```

PRP:NP → il # he
        VB → va # go
RB+VB → ne va pas # not go
        VP → ne va pas # does not go
        S → il ne va pas # he does not go

```

where the lhs symbol is separated from the source and target words by the \rightarrow symbol. Based on these initial rules, we perform the rule *generalization* procedure from (Chiang, 2005), replicated below: For each rule:

$$N \rightarrow f_1 \dots f_m \# e_1 \dots e_n$$

for which is an *initial* rule

$$M \rightarrow f_i \dots f_u \# e_j \dots e_v$$

where $1 \leq i < u \leq m$ and $1 \leq j < v \leq n$, a new rule can be generated that has the form:

$$N \rightarrow f_1 \dots f_{i-1} M_k f_{u+1} \dots f_m \# e_1 \dots e_{j-1} M_k e_{v+1} \dots e_n$$

where k is an index for the nonterminal M that indicates the one-to-one correspondence between the new M tokens on the two sides (it is not in the space of word indices like i, j, u, v, m, n). This generalization procedure can be performed recursively to create rules with multiple nonterminal symbols. The following rules would be extracted from our exam-

ple sentence:

$$\begin{aligned}
 S &\rightarrow \text{PRP:NP}_1 \text{ ne va pas } \# \text{ PRP:NP}_1 \text{ does not go} \\
 S &\rightarrow \text{il ne VB}_1 \text{ pas } \# \text{ he does not VB}_1 \\
 S &\rightarrow \text{il VP}_1 \# \text{ he VP}_1 \\
 S &\rightarrow \text{il RB+VB}_1 \# \text{ he does RB+VB}_1 \\
 S &\rightarrow \text{PRP:NP}_1 \text{ VP}_2 \# \text{ PRP:NP}_1 \text{ VP}_2 \\
 S &\rightarrow \text{PRP:NP}_1 \text{ RB+VB}_2 \# \text{ PRP:NP}_1 \text{ does RB+VB}_2 \\
 \text{VP} &\rightarrow \text{ne VB}_1 \text{ pas } \# \text{ does not VB}_1 \\
 \text{RB+VB} &\rightarrow \text{ne VB}_1 \text{ pas } \# \text{ not VB}_1 \\
 \text{VP} &\rightarrow \text{RB+VB}_1 \# \text{ does RB+VB}_1
 \end{aligned}$$

Nonterminal labels provide strong syntactic constraints in the grammar, but the heuristic non-terminal labels introduce significant sparsity. The importance of these constraints might vary across language pairs; we would like our model to determine the relative importance of these labels. Towards accomplishing this goal, for every labeled rule the SAMT grammar, we can also generate a non-syntactic rule labeled only with generic X nonterminals, like those in (Chiang, 2005). We introduce an additional feature in the log-linear translation model that allows the decoder to prefer labeled or unlabeled derivations. To suppress the creation of generic rules, pass the parameter ‘-generate_generic_variant 0’.

The number of rules generated by this procedure is exponential in the number of initial phrases pairs, producing a grammar that is impractical for efficient translation. The following parameters are used to restrict the number of rules extracted per sentence:

- *-max_abstraction_count* (default: 2): maximum number of abstractions (nonterminal pairs) per rule.
- *-max_source_symbol_count* (default: 6): maximum number of symbols (terminals and nonterminals) on the source side of the rule.

This restricted rule set can be pruned further with the following parameters for **MergeRules**:

- *-allow_consec_nts* (default: 1): if set to 0, discards rules that have consecutive nonterminals on the source side.
- *-allow_src_abstract* (default: 1): if 0, discards rules that do not have any source terminal symbols for example: $S \rightarrow \text{NP}_1 \text{ VP}_2 \# \text{ NP}_2 \text{ VP}_1$. Setting this parameter to 0, drastically reduces decoding time.
- *-nonlexminfreq*, *-lexminfreq* (defaults: 0): minimum occurrence frequency thresholds for non-lexical and lexical rules respectively. Increasing these thresholds reduces the size of the grammar, but often at the cost of translation quality (Zollmann et al., 2008).
- *-min_freq_given_src_arg* (default: 0): minimum relative frequency of a rule given its labeled source.

The labeling and extraction procedures defined above identify rules from the input word-aligned parallel corpora and associated parse trees. The occurrence counts from this extraction

process are used in estimating relative frequency features for each rule. The estimation of these features is described in the next section.

2.2. PSCFG Features

Given a source sentence f and a PSCFG grammar, the translation task can be expressed analogously to monolingual parsing with a CFG. We find the most likely derivation D of the input source sentence and read off the English translation, identified by composing α from each rule used in the derivation. This search for the most likely derivation can be defined as:

$$\hat{e} = \text{tgt} \left(\underset{D \in \text{Derive}(G): \text{src}(D)=f}{\text{arg max}} p(D) \right) \quad (1)$$

where $\text{tgt}(D)$ refers to the sequence of target terminal symbols generated by the derivation D , $\text{src}(D)$ refers to the source terminal symbols of D and $\text{Derive}(G)$ is the set of sentence spanning derivations of grammar G . The distribution p over derivations is defined by a log-linear model. The probability of a derivation D is defined in terms of the rules r that are used in D :

$$p(D) = \frac{p_{\text{LM}}(\text{tgt}(D))^{\theta_{\text{LM}}} \prod_{r \in D} \prod_{i=1}^m \lambda_i(r)^{\theta_i}}{Z(\theta_{\text{LM}}, \theta_1, \dots, \theta_m)} \quad (2)$$

where $\lambda_i(r)$ refers to features defined on each rule, p_{LM} is an n-gram language model (LM) probability distribution over target word sequences, and Z is a normalization constant that does not need to be computed during search under the arg max search criterion in Equation 1. The feature weights $\theta_{\text{LM}}, \theta_1, \dots, \theta_m$ are trained in concert with the language model weight via Minimum Error Rate (MER) training (Och, 2003). The features $\lambda_i(r)$ are statistics estimated from rule occurrence counts. They represent multiple criteria by which the decoder can judge the quality of each rule and, by extension, each derivation.

The Reduce step (**MergeRules** binary) of the Rule Extraction phase is responsible for generating the following features λ_i :

- $p(r | \text{lhs}(r))$: Probability of a rule given its lhs label.
- $p(r | \text{src}(r))$: Probability of a rule given its source side.
- $p(\text{ul}(\text{tgt}(r)) | \text{ul}(\text{src}(r)))$: Probability of the un-labeled target side of the rule given its un-labeled source side.

where lhs returns the left-hand-side of a rule, src returns the source side γ , tgt returns the target side α , and ul removes all *labels* from nonterminal symbols. For example, $\text{ul}(\text{NP}+\text{AUX}_1 \text{ does not go}) = \text{X}_1 \text{ does not go}$. The feature $p(\text{ul}(\text{tgt}(r)) | \text{ul}(\text{src}(r)))$ is equivalent to the target-given-source relative frequency estimate commonly used in phrase based systems for purely lexical rules, and the ul function allows us to calculate these estimates for rules with labeled nonterminal symbols as well.

To estimate the features above, we use maximum likelihood estimation based on counts of the rules extracted from the training data. For example, $p(r | \text{lhs}(r))$ is estimated by computing $\text{cnt}(r) / \text{cnt}(\text{lhs}(r))$, aggregating counts from all extracted rules. These relative frequency

features can be efficiently calculated by taking advantage of the MapReduce key sorting mechanism.

The output of the Rule Extraction phase is a grammar with a small set of features that has been learned automatically from the input data. The resulting grammar is large, and for most translations tasks, cannot be loaded directly into memory for decoding. To avoid this problem, the SAMT toolkit filters the grammar against a specific test corpus, generating a sentence specific grammar for each *sentence* in the corpus. This filtering is performed for each corpora that we need for translation, typically *development*, *test*, and *unseen test* corpora are used to train and evaluate machine translation systems.

2.3. Rule and LM Filtering

The Rule Filtering phase (binaries **MapSubsampleRules**, **filterrules_bin**) take as input: the grammar from the Rule Extraction phase, a corpus to filter the grammar against, and additional model files (such as translation lexica) to generate additional rule features λ_i . In the Map step, the grammar is filtered on a per-sentence basis by matching the source words of each rule to the source words in the sentence we want to translate. In the Reduce step, rules are augmented with the following features:

- $\lambda_{lex}(r) = p_w(\text{src}(r)|\text{tgt}(r)), p_w(\text{tgt}(r)|\text{src}(r))$: lexical weights based on terminal symbols as in (Koehn, Och, and Marcu, 2003).
- $\lambda_{glue}(r) = 1$ if the rule is a Glue rule as in (Chiang, 2005) that serves to monotonically concatenate span translations, 0 otherwise. The Glue rule is manually added to the grammar as described below.
- $\lambda_{ra}(r) = 1$ for all rules. A rule application count, allowing the model to favor derivations with more or less rules depending on the weight assigned to this feature.
- $\lambda_{tgt}(r)$: Count of the number of target terminals in r . Allows the model to prefer longer or shorter translations.
- $\lambda_{lex}(r) = 1$ if the rule has no nonterminals, 0 otherwise.
- $\lambda_{abs}(r) = 1$ if the rule has no terminals, 0 otherwise.
- $\lambda_{adj}(r) = 1$ if the rule has adjacent nonterminals in γ , 0 otherwise. Allows the model to indicate confidence in derivations that include multiple sequential nonterminals.
- $\lambda_X(r) = 1$ if the rule's lhs label is X (the hierarchical 'backoff' label), 0 otherwise
- $\lambda_{bal}(r) = 1$ if the ratio of source vs. target terminals in r is significantly different from the same ratio measured over sentences in the corpus, 0 otherwise.
- $\lambda_{mono}(r) = 1$ if the rule does not re-order its nonterminals, 0 otherwise.
- $\lambda_{rare}(r) = e^{(1/\text{cnt}(r))}$: uses the number of times a rule has been seen during training, $\text{cnt}(r)$, to allow penalization of derivations that use rare rules.

The Reduce step of Rule Filtering provides several options to further restrict the grammar and to augment the additional features. These options can be specified via the top-level parameter: *filter_params*. Documentation regarding these additional options can be found in the script: **filter_rules.pl** which is used to generate the MapReduce binary **filter_rules_bin**.

The Rule Filtering Reduce step also adds the following system rules to each sentence specific grammar.

- Beginning-of-sentence rule: $S \rightarrow \langle s \rangle \# \langle s \rangle$
- Glue rules (Chiang, 2005) for each NT N in the grammar, for example: $S \rightarrow S_1 N_2 \# S_1 N_2$
- End-of-sentence rule: $S \rightarrow S_1 \langle \backslash s \rangle \# S_1 \langle \backslash s \rangle$
- ‘Unknown’-rules (e.g. $NNP \rightarrow _UNKNOWN \# _UNKNOWN$) generating a limited set of labels for the word ‘_UNKNOWN’, which the decoder substitutes for unknown source words

The Glue rules (Chiang, 2005) play an important role in grammar based approaches to MT. These rules serve to simply concatenate translations of consecutive spans during decoding, similar to monotone decoding in a phrase based system (Koehn, Och, and Marcu, 2003). These Glue operations allow the system to produce translations that violate the syntactic constraints encoded in the labels of the grammar—at a cost determined via the MER trained weight θ_{glue} .

Building sentence specific grammars allows us to estimate the parameters and features of the grammar on large parallel corpora, while still being able to load all relevant rules to translate particular sentences in a test corpus. We follow this same approach to filter large n-gram language models in a LM Filtering phase. While the Rule Filtering phase filters rules based on the source side of the rule, the n-gram LM must be filtered according to the possible set of target words that can generated by applying the sentence specific grammar. For each sentence specific grammar, a possible target vocabulary is generated, which is used by the Rule Filtering binary (**LMFilter**) to produce sentence specific language models.

3. PSCFG Decoding

The sentence specific grammars and language models built via the MapReduce phases described above are used in a bottom-up chart parsing decoder to perform the search in Equation 1. The SAMT toolkit provides an implementation of the CYK+ algorithm (Chappelier and Rajman, 1998), that allows efficient decoding for grammars with more than two non-terminal symbols. Our decoder integrates n-gram language models during search, using the Cube Pruning algorithm described in (Chiang, 2007) to mitigate the computational impact of this feature. Decoding is performed as a MapReduce phase as well; taking as input individual sentences from the corpora that we want to translate, and translating them in the Map step (using the **FastTranslateChart** binary). Each decoder task has access to the sentence specific models that were built in previous phases on HDFS. The output of the Map step is a n-best list of candidate translations for the input source sentence. This n-best list is used in MER training. In order to facilitate MER training, the Reduce step in this phase merges multiple n-best lists across iterations for the same sentence—as required by MER training. The Map step takes several parameters that govern translation runtime and translation quality.

The runtime complexity of our decoder’s search is:

$$\mathcal{O} \left(|f|^3 \left[|\mathcal{N}| |\mathcal{T}_T|^{2(n-1)} \right]^K \right) \quad (3)$$

where K is the maximum number of NT symbols per rule, $|f|$ is the source sentence length, \mathcal{N} is the set of nonterminal labels in the grammar, \mathcal{T}_T is the set of target language terminals in the grammar, and n is the order of the n -gram LM. The grammar restriction parameters described in Section 2.1 have a large significant impact on this runtime (particularly *allow_src_abstract*, *allow_consec_nts*, *max_abstraction_count*), but this search still requires additional pruning to produce translations in reasonable time-frames—especially when translating longer sentences. The most important decoder parameters are described below:

- *wts*: corresponds to the weights θ in the translation model in Equation 2. In practice, these weights are iteratively trained via MER.
- *HistoryLength*: (default 2) The number of words considered as LM history length during decoding. When set to less than $n - 1$, when using an n -gram LM, decoding time is reduced at the expense of search errors, which can reduce translation quality.
- *SRIHistoryLength*: This value indicates the full history length of the n -gram language model. When using a reduced *HistoryLength*, this value is used to recover from search errors in a LM-driven n -best extraction step similar to (Huang and Chiang, 2007).
- *PruningMap*: (default: 0-100-5-@_S-200-5): Format: *lhs-b-β*. Pruning parameters for Cube Pruning (Chiang, 2007). For each nonterminal label *lhs* in the grammar for a source span during decoding, this parameter restricts the number of chart items to b items, and items that are have cost of at most β greater than the best item. *lhs = 0* sets pruning parameters for all *lhs* symbols that have not been explicitly specified.
- *ComboPruningBeamSize* : (default 10000) Sets the maximum number of items generated in each cell via Cube Pruning. Reducing this value reduces decoding time when *PruningMap* limits have not caused pruning.
- *MaxHypsPerCell*: (default 1000000000) Limits the total number of items (partial translation hypotheses) created for each span during decoding—across items that have different *lhs* labels (not counting X and S items, which always pass thru this pruning filter). This value is typically set when using grammars with a large number of *lhs* labels to reduce translation runtime, but does introduce additional search error.
- *MaxCostDifferencePerCell*: (default inf) Max. allowed cost that an item can deviate from the best item in its chart cell (inf: any cost allowed). Items with *lhs* X or S always pass thru this filter. This and the previous parameter are the only parameters that apply pruning across items with different nonterminal labels.
- *MaxCombinationCount*: (default 10) Limits the application of automatically learned PSCFG rules to source spans less than or equal to *MaxCombinationCount*. Spans of greater length are composed monotonically with Glue rules. Decoding time is linear in sentence length once this limit is in effect.

Track	Words (English)	LM 1-N grams (N)	Dev.	Test1	Test2
IWSLT	632K	431,292 (5)	IWSLT06	IWSLT07	N/A
67M	67M	102,924,025 (4)	MT05	MT06	MT08
230M	230M	273,233,010 (5)	MT05	MT06	MT08

Table 1: Training data configurations used to evaluate SAMT on Hadoop. The number of words in the target text and the number of 1-N grams represented in the complete model are the defining statistics that characterize the scale of each task. For each LM we also indicate the order of the n-gram model.

3.1. Minimum Error Rate Training

The parameters θ are trained via Minimum Error Rate (MER) training (Och, 2003) to maximize translation quality according to a user specified automatic translation metric, like BLEU (Papineni et al., 2002) or NIST (Doddington, 2002). MER training is implemented in the SAMT toolkit as a MapReduce phase using n-best lists from the decoding phase. The decoder annotates each candidate translation for a given source sentence with statistics that allow the MER procedure to evaluate evaluation metric error surfaces. The SAMT implementation of MER (binary: *MER*) also performs feature selection, where sparse solutions ($\theta_i = 0$) to the optimization procedure are preferred. The following parameters initiate and affect MER training (to get information about additional parameters, run *MER* without any input).

- *mer*: (default false) A top-level parameter to initiate MER training on a development corpus.
- *ScoringMetric*: (default IBMBLEU) The automatic evaluation metric that MER optimizes towards.
- *DisplayNBestData*: (default 1000) The size of the n-best list used for MER training.
- *Opti_Epsilon*: (default 0.0001) When automatic metric scores differ by less than this parameter, MER training is terminated. This is also the significance threshold when testing for sparse solutions.

4. Empirical Results

We demonstrate the SAMT on Hadoop toolkit on three Chinese-to-English translation tasks, representing a wide range of resource conditions. Each task is described in Table 1. The IWSLT task is a limited resource, limited domain task, while 67M and 230M (named for their respective corpora sizes), are corpora used for the annual NIST MT evaluation. For each task we list the number of words in the target side of the corpus and the number of 1-n grams in the n-gram LM (estimated from parallel and monolingual data).

For each resource condition, we build SAMT systems using a purely hierarchical grammar (Hier) (Chiang, 2005) and a syntax augmented grammar (Syntax) from (Zollmann and Venugopal, 2006). All experiments use a 2-gram *HistoryLength* length the first pass of decoding, and the full LM history during the second pass n-best list search. These grammars are built with ‘*-allow_consec_nots 0 -allow_src_abstract 0*’, and the NIST MT task rules are addition-

System	Dev. BLEU	Test1 BLEU	Test2 BLEU	Grammar Train. (h:m)	Test1 (m)
IWSLT Hier	27.0	37.0	N/A	0:12	4
IWSLT Syntax	30.9	37.2	N/A	0:26	12
67M Hier	35.19	32.98	25.88	1:10	17
67M Syntax	35.69	33.12	26.48	2:26	65
230M Hier	36.39	33.74	26.28	4:13	23
230M Syntax	37.11	34.04	26.74	7:21	53

Table 2: Translation quality as measured by IBM-BLEU% (i.e., brevity penalty based on closest reference length) on each resource track for appropriate evaluation data sets. Systems 67M and 230M are evaluated in lower-case, while IWSLT is evaluated in mixed case. Training and decoding times given are based on a cluster of 100 1.9GHz Intel Xeon processors.

ally restricted by `'-nonlexminfreq 2 -min_freq_given_src_arg α '` where $\alpha = 0.005$ (Hier) and $\alpha = 0.01$ (Syntax). The Syntax based systems also use `'-MaxHypsPerCell 1000'` to limit the run time impact of the large number of lhs labels in these grammars.

In Table 2, we report BLEU scores on development and test data as well as run times to train the respective PSCFG grammars and perform translation with them. Training run times are reported based on Hadoop MapReduce jobs running on a cluster of 50 dedicated machines, each running 2 Map or Reduce tasks each. These results demonstrate the ability for the SAMT toolkit to scale to large resource data conditions. For each of the the three data conditions we see that training the Syntax grammar takes longer to train as well as translate with. Translation quality improvements that result from using more parallel and monolingual data are clear when comparing the 67M and 230M systems. In these experiments, we see small but consistent improvements from the introduction of SAMT labels, in line with experiments in (Zollmann et al., 2008). Overall, translation quality results reported here are competitive with reported results in the literature and constitute a valid baseline for further research.

5. Conclusions and Resources

In this paper we have described the SAMT on Hadoop toolkit, an end-to-end framework for large scale grammar based statistical machine translation. We discussed the methodology of the SAMT approach, and described important toolkit parameters that affect translation quality and run time. Built upon the open-source Hadoop distributed computation framework, our toolkit is able to scale to build grammars for large scale translation tasks in reasonable time frames. The toolkit can be easily extended to experiment with alternative grammar extraction and decoding techniques.

Additional documentation to download and build SAMT on Hadoop is available on the toolkit website. Example parameter files for the IWSLT task are included in the distribution. SAMT on Hadoop assumes a fully functional Hadoop installation, using Hadoop-on-Demand to allocate clusters of dedicated machine for each computational task.

Bibliography

- Brown, Peter F., Vincent J. Della Pietra, Stephen A. Della Pietra, and Robert L. Mercer. 1993. The mathematics of statistical machine translation: parameter estimation. *Computational Linguistics*.
- Chappelier, J.C. and M. Rajman. 1998. A generalized CYK algorithm for parsing stochastic CFG. In *Proceedings of Tabulation in Parsing and Deduction (TAPD)*, pages 133–137, Paris.
- Charniak, Eugene. 2000. A maximum entropy-inspired parser. In *Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics Conference (HLT/NAACL)*.
- Chiang, David. 2005. A hierarchical phrase-based model for statistical machine translation. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*.
- Chiang, David. 2007. Hierarchical phrase based translation. *Computational Linguistics*.
- Chiang, David and Kevin Knight. 2006. An introduction to synchronous grammars. In *Tutorials at the Annual Meeting of the Association for Computational Linguistics (ACL)*.
- Cutting, Doug and Eric Baldeschwieler. 2007. Meet Hadoop. In *O'Reilly Open Software Convention*, Portland, OR.
- Dean, Jeffrey and Sanjay Ghemawat. 2004. Mapreduce: Simplified data process on large cluster. In *Proceedings of Symposium on Operating System Design and Implementation*.
- Doddington, George. 2002. Automatic evaluation of machine translation quality using n-gram co-occurrence statistics. In *In Proceedings ARPA Workshop on Human Language Technology*.
- Dyer, Christopher, Aaron Cordova, Alex Mont, and Jimmy Lin. 2008. Fast, easy, and cheap: Construction of statistical machine translation models with mapreduce. In *Proceedings of the Workshop on Statistical Machine Translation, ACL*.
- Galley, Michael, Mark Hopkins, Kevin Knight, and Daniel Marcu. 2006. Scalable inferences and training of context-rich syntax translation models. In *Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics Conference (HLT/NAACL)*.
- Huang, Liang and David Chiang. 2007. Forest rescoring: Faster decoding with integrated language models. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*.
- Koehn, Philipp, Franz J. Och, and Daniel Marcu. 2003. Statistical phrase-based translation. In *Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics Conference (HLT/NAACL)*.
- Och, Franz J. 2003. Minimum error rate training in statistical machine translation. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*.
- Och, Franz J. and Hermann Ney. 2003. A systematic comparison of various alignment models. *Computational Linguistics*.
- Och, Franz J. and Hermann Ney. 2004. The alignment template approach to statistical machine translation. *Computational Linguistics*.

- Papineni, Kishore, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: a method for automatic evaluation of machine translation. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*.
- Paul, Michael. 2006. Overview of the IWSLT 2006 evaluation campaign. In *Proceedings of the International Workshop on Spoken Language Translation (IWSLT)*.
- Petrov, Slav, Aria Haghighi, and Dan Klein. 2008. Coarse-to-fine syntactic machine translation using language projections. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Steedman, Mark. 1999. Alternative quantifier scope in CCG. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*.
- Stolcke, Andreas. 2002. SRILM —an extensible language modeling toolkit. In *Proceedings of the International Conference on Spoken Language Processing (ICSLP)*.
- Wu, Dekai. 1997. Stochastic inversion transduction grammars and bilingual parsing of parallel corpora. *Computational Linguistics*.
- Zhang, Hao and Daniel Gildea. 2008. Efficient multi-pass decoding for synchronous context free grammars. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*.
- Zollmann, Andreas and Ashish Venugopal. 2006. Syntax augmented machine translation via chart parsing. In *Proceedings of the Workshop on Statistical Machine Translation, HLT/NAACL*, New York, June.
- Zollmann, Andreas, Ashish Venugopal, Franz J. Och, and Jay Ponte. 2008. A systematic comparison of phrase-based, hierarchical and syntax-augmented statistical MT. In *Proceedings of the Conference on Computational Linguistics (COLING)*.
- Zollmann, Andreas, Ashish Venugopal, and Stephan Vogel. 2008. The CMU Syntax-Augmented Machine Translation System: SAMT on Hadoop with N-best Alignments. In *Proc. of the International Workshop on Spoken Language Translation*, pages 18–25, Hawaii, USA.