
Fast, Scalable Phrase-Based SMT Decoding

Abstract

The utilization of statistical machine translation (SMT) has grown enormously over the last decade, many using open-source software developed by the NLP community. As commercial use has increased, there is need for software that is optimized for commercial requirements, in particular, fast phrase-based decoding and more efficient utilization of modern multicore servers.

In this paper we re-examine the major components of phrase-based decoding and decoder implementation with particular emphasis on speed and scalability on multicore machines. The result is a drop-in replacement for the Moses decoder which is up to fifteen times faster and scales monotonically with the number of cores.

1 Introduction

SMT has steadily progressed from a research discipline to commercial viability during the past decade as can be seen from services such as the Google and Microsoft Translation services. As well as general purpose services such as these, there is a large number of companies that offer customized translation systems, as well as companies and organization that implement in-house solutions. Many of these customized solutions use Moses as their SMT engine.

For many users, decoding is the most time-critical part of the translation process. Making use of the multiple cores that are now ubiquitous in today's servers is a common strategy to ameliorate this issue. However, it has been noticed that the Moses decoder, amongst others, is unable to efficiently use multiple cores (Fernández et al., 2016). That is, decoding speed does not substantially increase when more cores are used, in fact, it may actually *decrease* when using more cores. There have been speculation on the causes of the inefficiency as well as potential remedies.

This paper is the first we know of that focuses on improving decoding speed on multicore servers. We will take a holistic approach to solving this issue. We create a decoder that is optimized for multi-core processing speed by concentrating on four main areas.

1. Faster memory management of data-structures through the use of customized memory pools
2. Exploring alternatives to cardinality-based hypothesis stack configuration
3. Re-examining the efficiency of phrase-table lookup using translation rule caching and data compression
4. Integrating the lexicalized re-ordering model into the phrase-table, thus eliminating the need for independent random lookup this model

The result is a decoder that is significantly faster than the Moses baseline for single-threaded operation, and scales with the number of cores.

As far as possible, model scores and functionality are compatible with Moses to aid comparison and ease transition for existing users. All source code will be made available under an open-source license.

The rest of the paper will be broken up into the following sections. The rest of this section will discuss prior work and an outline of the phrase-based model. Section 2 will then describe the modifications to improve decoding speed. We describe the experiment setup in Section 3 and present results Section 4. We conclude in the last section and discuss possible future work.

1.1 Prior Work

Most prior work on increasing decoding speed look to optimizing specific components of the decoder or the decoding algorithm.

Heafield (2011) and Yasuhara et al. (2013) describes fast, efficient datastructures for language models. Zens and Ney (2007) describes an implementation of a phrase-table for an SMT decoder that is loaded on demand, reducing the initial loading time and memory requirements. Junczys-Dowmunt (2012) extends this by compressing the on-disk phrase table and lexicalized re-ordering model.

Chiang (2007) describes the cube-pruning and cube-growing algorithm which allows the tradeoff between speed and translation quality to be adjusted with a single parameter. Wuebker et al. (2012b) note that language model querying is amongst the most expensive operation in decoding. They sought to improved decoding speed by caching score computations early pruning of translation options. This work is similar to Heafield et al. (2014) which group hypotheses with identical language model context and incrementally expand them, reducing LM querying.

Fernández et al. (2016) was concerned with multi-core speed but treated the decoding process as a black box within a parallelization framework.

There are a number of phrase-based decoding implementations, many of which implements the extensions to the phrase-based model described above.

The most well known is Moses (Koehn et al., 2007) which implements a number of speed optimizations, including cube-pruning. It is widely used for MT research and commercial use.

Joshua (Li et al., 2009) also supports cube-pruning for phrase-based models. Phrasal (Spence Green and Manning, 2014) supports a number of variants of the phrase-based model. Jane (Wuebker et al., 2012a) supports the language model look-ahead described in Wuebker et al. (2012b) in addition to other tools to speed up decoding such as having a separate fast, lightweight decoder. mtplz is a specialized decoder developed to implement the incremental decoding described in Heafield et al. (2014).

The Moses, Joshua and Phrasal decoders implement multithreading, however, they all report scalability problems, either in the paper (Phrasal) or via social media (Moses¹ and Joshua²).

Jane and mtplz are single-threaded decoders, relying on external applications to parallelize operations.

This paper not only focuses on faster single-threaded decoding but also on overcoming the shortcomings of existing decoding implementations on multicore servers.

1.2 Phrase-Based Model

The objective of decoding is to find the target translation with the maximum probability, given a source sentence. That is, for a source sentence s , the objective is to find a target translation \hat{t} which has the highest conditional probability $p(t|s)$. Formally, this is written as:

$$\hat{t} = \arg \max_t p(t|s) \quad (1)$$

where the *arg max* function is the search. The log-linear model generalizes Equation 1 to include more component models and weighting each model according to the contribution of

¹<https://github.com/moses-smt/mosesdecoder/issues/39>

²<https://twitter.com/ApacheJoshua/status/342022794097340416>

each model to the total probability.

$$p(t|s) = \frac{1}{Z} \exp\left(\sum_m \lambda_m h_m(t, s)\right) \quad (2)$$

where λ_m is the weight, and h_m is the *feature function*, or ‘score’, for model m . Z is the partition function which can be ignored for optimization.

The standard feature functions in the phrase-based model include:

1. a distortion penalty
2. a phrase-penalty,
3. a word penalty,
4. an unknown word penalty.
5. log transforms of the target language model probability $p(t)$,
6. log transforms translation model probabilities, $p_{TM}(t|s)$ and $p_{TM}(s|t)$, and word-based translation probabilities $p_w(t|s)$ and $p_w(s|t)$,
7. log transforms of the lexicalized re-ordering probabilities,

Of these feature functions, we will focus on optimizing the speed of the phrase-table and lexicalized re-ordering models.

1.3 Beam Search

A translation of a source sentence is created by applying a series of translation rules which together translate each source word once, and only once. Each partial translation is known as a *hypothesis*, which is created by applying a rule to an existing hypothesis. This *hypothesis expansion* process starts with a hypothesis that has translated no source word and ends with completed hypotheses that has translated all source words. The highest-scoring completed hypothesis, according to the model score, is considered the best translation, \hat{t} .

In the phrase-based model, each rule translates a contiguous sequence of source words. Successive applications of translation rules do not have to be adjacent on the source side, depending on the distortion limit. The target output is constructed strictly left-to-right from the target side from the series of translation rules.

A beam search algorithm is used to create the completed hypothesis set efficiently. Hypotheses are grouped into stacks where each stack holds a number of comparable hypotheses. Most phrase-based implementations group hypotheses according to coverage cardinality.

2 Proposed Improvements

We will also concentrate on four main areas for optimization.

2.1 Efficient Memory Allocation

The search algorithm creates and destroy a large number of intermediate objects such as hypotheses and feature function states. The memory management this requires puts a burden on the operating system due to the need to synchronize memory access. This problem worsens with multi-threaded applications such as Moses.

Libraries such as tcmalloc (Ghemawat and Menage, 2009) are designed to reduce locking contention for multi-threaded application but in our case, this is still not enough.

We shall seek to improve decoding speed by replacing the operating system’s general purpose memory management with our own custom memory management scheme. Memory will be allocated from a memory pool rather than use the operating system’s general purpose allocation functions.

A memory pool is a large block of memory that has been given to the application by the operating system. The application is then responsible for allocating portions of this memory to its components as required. We will use thread-specific memory pools to increase speed by avoiding locking contention during memory access. Our memory pools will be dynamic, growing in size when required but never reducing in size. The pools are not deleted until the decoder application ends.

To further increase memory management speed, objects in the memory pool are not deleted. Unused data structures accumulates in the pool until a reset event. After the event, the pool is assumed to be empty and simply reused.

We instantiate two memory pools per decoding thread, one which is never reset and another which is reset after the decoding of each sentence. Data structures are created in either pool according to their life cycle.

Accumulating unused objects in the memory pools can result in unacceptably high memory usage so object queues are available for high-churn objects which allows the decoder to re-cycle unused objects before the reset event. This also increase speed as LIFO queues are used so that the most recently accessed memory are used, increasing CPU cache hits.

2.2 Stack Configurations

The most popular stack configuration for phrase-based models is coverage cardinality, that is, hypotheses that have translated the same number of source words are stored in the same stack. This is implemented in Pharaoh, Moses and Joshua.

However, there are alternatives to this configuration. Och et al. (2001) uses a single stack for all hypotheses, Brown et al. (1993) uses coverage stacks (ie. one stack per unique coverage vector) while Wuebker et al. (2012a) and Zens and Ney (2008) apply both coverage and cardinality pruning. While useful, these prior works present only one particular stack configuration each.

Ortiz-Martínez et al. (2006) explore a range of stack configurations by defining a granularity parameter which controls the maximum number of stacks required to decode a sentence.

We shall re-visit the question of stack configuration with a particular emphasis on how they can help improve the tradeoff between speed and translation quality. We will do so in the context of the cube-pruning algorithm, the algorithm that we will be using and which was not available to many of the earlier work.

2.3 Phrase-Table Optimizations

For any phrase-table table of a realistic size, memory and loading time constraints requires us to use a load-on-demand implementation. Moses has several which we can make use of, each with differing performance characteristics. Figure 1 shows the decoding speed for the fastest two implementations. From this, it appears that the Probing phrase-table (Bogoychev and Lopez, 2016) has the fastest translation rule lookup, especially with large number of cores, therefore, we will concentrate exclusively on this implementation from hereon.

We propose two optimizations. Firstly, the translation rule caching mechanism in Moses saves the most recently used rules. However, this require locking and active management in clearing of old rules. The result is *slower* decoding, Table 1.

	No cache	Caching
Decoding time	2877	2540 (-12%)

Table 1: Decoding speed (in words / sec with 32 threads) when using phrase-table cache

We shall explore a simpler caching mechanism by creating a static cache of the most likely

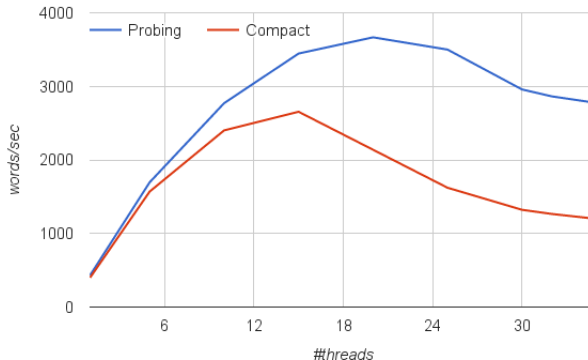


Figure 1: Moses decoding speed with two different phrase-table implementations

translation rules to be used at the start of decoding.

Secondly, the Probing phrase-table use a simple compression algorithm to compress the target side of the translation rule. Compression was championed by Junczys-Dowmunt (2012) as the main reason behind the speed of their phrase-table but as we saw in Figure 1, this comes at the cost of scalability to large number of threads. We shall therefore take the opposite approach to and improve decoding speed by disabling compression.

2.4 Lexicalized Re-ordering Model Optimizations

Similar to the phrase-table, the lexicalized re-ordering model is trained on parallel data. A resultant model file is then queried during decoding. The need for random lookup during querying inevitably results in slower decoding speed. Previous work such as Junczys-Dowmunt (2012) improve querying speed with more efficient data structures.

However, the model’s query keys are the source and target phrase of each translation rule. Rather than storing the lexicalized re-ordering model separately, we shall integrating it into the translation model, eliminating the need to query a separate file. However, the model remain the same under the log-linear framework, including having its own weights.

This optimization has precedent in Wuebker et al. (2012a) but the effect on decoding speed were not published. We will compare results with using a separate model in this paper.

3 Experimental Setup

We trained a phrase-based system using the Moses toolkit with standard settings. The training data consisted of most of the publicly available Arabic-English data from Opus (Tiedemann, 2012) containing over 69 million parallel sentences, and tuned on a held out set. The phrase-table was then pruned, keeping only the top 100 entries per source phrase, according to $p(t|s)$. All model files were then binarized; the language models were binarized using KenLM (Heafield, 2011), the phrase table using the Probing phrase-table, lexicalized re-ordering model using the compact data structure. These binary formats were chosen for their best-in-class multithreaded performance. Table 2 gives details of the resultant sizes of the model files. For testing decoding speed, we used a subset of the training data, Table 3.

For verification with a different dataset, we also used a second system trained on the French-English Europarl corpus (2m parallel sentences). The two different systems have characteristics that we are interested in analyzing; ar-en have short sentences with large models while fr-en have overly long sentences with smaller models. Where we need to compare model scores,

	ar-en	fr-en
Phrase table	17	5.8
Language model	3.1	1.8
Lex re. model	2.3	637MB

Table 2: Model sizes in GB

we used held out test sets.

	ar-en	fr-en
For speed testing		
Set name	Subset of training data	
# sentences	800k	200k
# words	5.8m	5.9m
Avg words/sent	7.3	29.7
For model score testing		
Set name	OpenSubtitles	newstest2011
# sentences	2000	3003
# words	14,620	86,162
Avg words/sent	7.3	28.7

Table 3: Test sets

Standard Moses phrase-based configurations are used, except that we use the cube-pruning algorithm (Chiang, 2007) with a pop-limit of 400³, rather than the basic phrase-based algorithm. The cube-pruning algorithm is often employed by users who require fast decoding as it gives them the ability to trade speed with translation quality via a simple pop-limit parameter.

As a baseline, we use a recent⁴ version of the Moses decoder taken from the github repository.

For all experiments, we used a Dell PowerEdge R620 server with 16 cores, 32 hyper-threads, split over 2 physical processors (Intel Xeon E5-2650 @ 2.00GHz). The server has 380GB RAM. The operating system was Ubuntu 14.04, the code was compiled with gcc 4.8.4 and Boost 1.59⁵ and the tcmalloc library.

4 Results

4.1 Optimizing Memory

Each sentence is decoded from start to finish using one worker thread.

As can be seen in Table 4, over 24% of the Moses decoder running time is spent on memory management and this increases to 39% when 32 threads are used, dampening the scalability of the decoder. By contrast, our decoder spends 11% on memory management and does not significantly increase with more threads.

Figure 2 compares the decoding speed for Moses and our decoder, using the same models, parameters and test set. Our decoder is over 3 times faster with one thread, and 5.0 times faster

³the pop-limit was chosen from public discussion on the Moses mailing list on an acceptable balance between decoding speed and translation quality with Moses for commercial use

⁴The experiments were performed between January and May 2016 with the latest github code to hand. The main ar-en experiments were rerun with the source code as of 8th June, 2016 to ensure there were no material difference. The commit hash was bc5f8d15c6ce4bc678ba992860bfd4be6719cee8

⁵<http://boost.org/>

	Moses		Our Work	
# threads	1	32	1	32
Memory	24%	39%	11%	13%
LM	12%	2%	47%	38%
Phrase-table	9%	5%	2%	4%
Lex RO	8%	2%	2%	2%
Search	2%	0%	14%	19%
Misc/Unknown	45%	39%	24%	29%

Table 4: Profile of %age decoding time

using all cores. However, like Moses, performance actually worsens after approximately 15 threads.



Figure 2: Decoding speed of Moses and our decoder, using the same models

4.2 Stack Configuration

We shall investigate the effects of the following three stack configurations on model score and decoding speed:

1. coverage cardinality,
2. coverage,
3. coverage and end position of most recently translated source word.

Coverage cardinality is the same as that in Moses and Joshua. Coverage configuration uses one stack per unique coverage vector.

Coverage and end position of most recently translated source word extends the coverage configuration by separating hypotheses where the last translate word are different. This is an optimization to reduce the number of checks on the distortion limit, which is dependent on the last word position.

The check is a binary function $d(C_h, e_{hypo}, range_r)$, where C_h is the coverage vector of hypothesis h , e_h is the end position of most recent source word that has been translated, and $range_r$ is the coverage of the rule to be applied. Figure 3 shows the standard hypothesis expansion algorithm.

By grouping hypotheses according to coverage *and* end position into groups we call '*min-istack*', the distortion limit only needs to be checked for each group, Figure 4.

```

for all  $h$  in  $stack_{|C|}$  do
  for all  $r$  in translation rules do
    if  $d(C_h, e_{hypo}, range_r)$  then
      expand  $h$  with  $r \rightarrow h'$ 
      add  $h'$  to stack  $stack_{C'}$ 
    end if
  end for
end for

```

Figure 3: Hypothesis Expansion with Cardinality Stacks

```

for all  $ministack_{C,e}$  in  $stack_{|C|}$  do
  for all  $r$  in translation rules do
    if  $d(C_h, e_{hypo}, range_r)$  then
      for all  $h$  in  $ministack_{C,e}$  do
        expand  $h$  with  $r \rightarrow h'$ 
        add  $h'$  to  $ministack_{C',e'}$ 
      end for
    end if
  end for
end for

```

Figure 4: Hypothesis Expansion with Coverage & End Position Stacks

Also, stack pruning occurs on each hypothesis group independently, therefore, potentially affect search errors and model scores.

Figure 5 present the tradeoff between decoding time and average model at various pop-limits. The model scores for all stack configurations are identical for low pop-limits but grouping hypotheses into coverage & end position produces higher model scores for higher pop-limits. It is also slower but the time/quality tradeoff is better overall with this stack configuration. Model scores for coverage and cardinality configurations are identical, which differs from the results in Ortiz-Martínez et al. (2006). This may be due to differing test set characteristics, or interactions with cube-pruning algorithm we used.

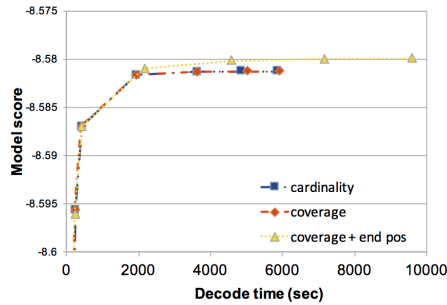


Figure 5: Trade-off between decoding time average model scores for different stack configurations

4.3 Translation Model

In the first optimization, we create a static translation model cache containing translation rules that translates the most common source phrases. This is constructed during phrase-table training based on the source counts. The cache is then loaded when the decoder is started. It does not require the overhead of managing an active cache but there is still some overhead in using a cache. Overall however, using a static cache result in a 10% decrease in decoding time if the

Cache size	Decoding Time	Cache Hit %age
Before caching	229	N/A
0	239 (+4.4%)	0%
1,000	213 (-7.0%)	11%
2,000	204 (-10.9%)	13%
4,000	205 (-10.5%)	14%
10,000	207 (-9.7%)	17%

Table 5: Decoding time (in secs with 32 threads) for varying cache sizes

optimum cache size is used, Table 5.

For the second optimization, we disable the compression of the target side of the translation rules. This increase the size of the binary files from 17GB to 23GB but the time saved not needing to decompress the data resulted in a 1.5% decrease in decoding time with 1 thread and nearly 7% when the CPUs are saturated, Table 6.

# threads	Compressed pt	Non-compressed pt
1	3052	3006 (-1.5%)
5	756	644 (-14.8%)
10	372	362 (-2.7%)
15	284	250 (-12.0%)
20	244	227 (-7.0%)
25	218	209 (-4.1%)
30	206	192 (-6.8%)
35	203	189 (-6.9%)

Table 6: Decoding time (in secs with 32 threads) for compressed and non-compressed phrase-tables

4.4 Lexicalized Re-ordering Model

The lexicalized re-ordering model requires a probability distribution of the re-ordering behaviour of each translation rule learnt from the training data. This is represented in the model file as a fixed number of probabilities for each rule, exactly how many probabilities is dependant on the model’s parameterization during training. During decoding, a probability from this distribution is assigned to each hypothesis according to the re-ordering of the translation rule.

Rather than holding the model probability distributions in the separate file, we pre-process the translation model file to include the lexicalized re-ordering model distributions for each rule. During decoding, the probability distribution is then taken from the translation model instead of querying a separate file.

This resulted in a significant decrease in decoding time, especially with high number of cores, Figure 6. Decoding time decreased by 29% with running with a single thread but it is over 5 times faster using 32 threads.

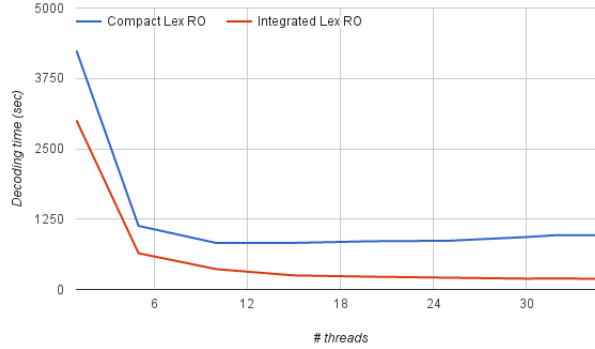


Figure 6: Decoding time with Compact Lexicalized Re-ordering, and integrated into a model the phrase-table

4.5 Scalability

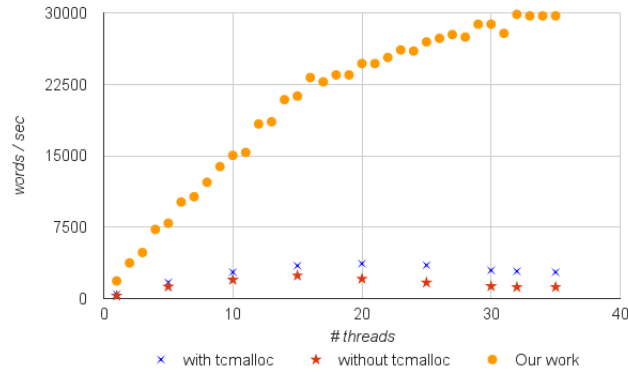


Figure 7: Comparison of decoding speed of our work and Moses (with & without the tcmalloc library)

Figure 7 shows decoding speed against the number of threads used. In our work, there is a constant increase in decoding speed when more threads are used, decreasing slightly after 16 threads when virtual cores are employed by the CPU. Overall, decoding is 16 times faster than single-threaded decoding when all 16 cores (32 hyperthreads) are fully utilized.

This contrast with Moses where speed increases to approximately 16 threads but then actually become slower thereafter. Using the tcmalloc library has a small positive effect on decoding speed but does little to improve scalability

Our work is 4.3 times faster than Moses with a single-thread and 10.4 faster when all cores are used.

5 Other Models and Even More Cores

Our decoder show no scalability issues when we tested with the same model and tested set on a larger server, Figure 8.

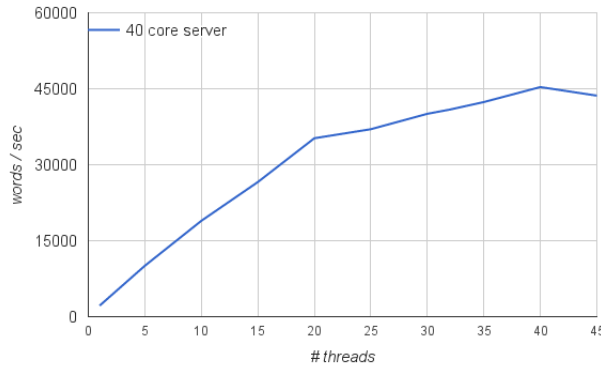


Figure 8: Decoding speed for with bigger servers

We verify the results with the French-English phrase-based system and test set. The speed gains are even greater than the Arabic-English test scenario, Figure 9. Our decoder is 5.4 times faster than Moses with a single-thread and 14.5 faster when all cores are saturated.

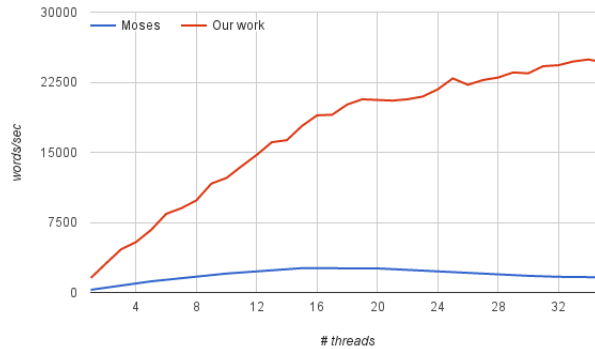


Figure 9: Decoding speed for fr-en model

It has been suggested that using a larger language model would overpower the improvements in decoding speed. We tested this conjecture by replacing the language model in the ar-en experiment with a 96GB language model. The time to load of language model is significant (394 sec) and was excluded from the translation speed. Results show that our decoder is 7 times faster than Moses and still scales monotonically until all CPUs are saturated, Figure 10.

6 Conclusion

We have presented a new decoder that is compatible with Moses. By studying the shortcomings of the current implementation, we are able to optimize for speed, particularly for multicore operation. This resulted in double digit gains compared to Moses on the same hardware. Our implementation is also unaffected by scalability issues that has afflicted Moses.

In future, we shall investigate other major components of the decoding algorithm, particularly the language model which has not been touched in this paper. We are also keen to explore the underlying reasons for the scalability issues in Moses to get a better understanding where

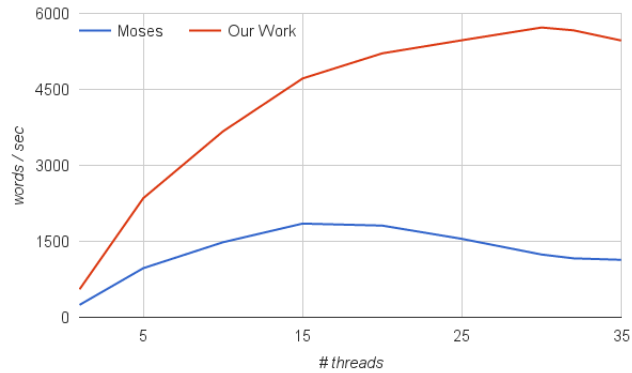


Figure 10: Decoding speed with a large language model

potential performance issues can arise.

References

- Bogoychev, N. and Lopez, A. (2016). Fast and highly parallelizable phrase table for statistical machine translation. In *Proceedings of the First Conference on Statistical Machine Translation WMT16*, Berlin, Germany.
- Brown, P. F., Della-Pietra, S. A., Della-Pietra, V. J., and Mercer, R. L. (1993). The mathematics of statistical machine translation. *Computational Linguistics*, 19(2):263–313.
- Chiang, D. (2007). Hierarchical phrase-based translation. *Computational Linguistics*, 33(2):201–228.
- Fernández, M., Pichel, J. C., Cabaleiro, J. C., and Pena, T. F. (2016). Boosting performance of a statistical machine translation system using dynamic parallelism. *Journal of Computational Science*, 13:37–48.
- Ghemawat, S. and Menage, P. (2009). Tcmalloc: Thread-caching malloc.
- Heafield, K. (2011). KenLM: faster and smaller language model queries. In *Proceedings of the EMNLP 2011 Sixth Workshop on Statistical Machine Translation*, pages 187–197, Edinburgh, Scotland, United Kingdom.
- Heafield, K., Kayser, M., and Manning, C. D. (2014). Faster Phrase-Based decoding by refining feature state. In *Proceedings of the Association for Computational Linguistics*, Baltimore, MD, USA.
- Junczys-Dowmunt, M. (2012). A space-efficient phrase table implementation using minimal perfect hash functions. In Sojka, P., Hork, A., Kopecek, I., and Pala, K., editors, *15th International Conference on Text, Speech and Dialogue (TSD)*, volume 7499 of *Lecture Notes in Computer Science*, pages 320–327. Springer.
- Koehn, P., Hoang, H., Birch, A., Callison-Burch, C., Federico, M., Bertoldi, N., Cowan, B., Shen, W., Moran, C., Zens, R., Dyer, C., Bojar, O., Constantin, A., and Herbst, E. (2007). Moses: Open source toolkit for statistical machine translation. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics Companion Volume Proceedings of the Demo and Poster Sessions*, pages 177–180, Prague, Czech Republic. Association for Computational Linguistics.

- Li, Z., Callison-Burch, C., Dyer, C., Khudanpur, S., Schwartz, L., Thornton, W., Weese, J., and Zaidan, O. (2009). Joshua: An open source toolkit for parsing-based machine translation. In *Proceedings of the Fourth Workshop on Statistical Machine Translation*, pages 135–139, Athens, Greece. Association for Computational Linguistics.
- Och, F. J., Ueffing, N., and Ney, H. (2001). An efficient A* search algorithm for statistical machine translation. In *Workshop on Data-Driven Machine Translation at 39th Annual Meeting of the Association of Computational Linguistics (ACL)*.
- Ortiz-Martínez, D., García-Varea, I., and Casacuberta, F. (2006). Generalized stack decoding algorithms for statistical machine translation. In *Proceedings on the Workshop on Statistical Machine Translation*, pages 64–71, New York City. Association for Computational Linguistics.
- Spence Green, D. C. and Manning, C. D. (2014). Phrasal: A toolkit for new directions in statistical machine translation. In *Proceedings of the Ninth Workshop on Statistical Machine Translation*, pages 114–121. Citeseer.
- Tiedemann, J. (2012). Parallel data, tools and interfaces in opus. In *LREC*, pages 2214–2218.
- Wuebker, J., Huck, M., Peitz, S., Nuhn, M., Freitag, M., Peter, J.-T., Mansour, S., and Ney, H. (2012a). Jane 2: Open source phrase-based and hierarchical statistical machine translation. In *24th International Conference on Computational Linguistics*, page 483. Citeseer.
- Wuebker, J., Ney, H., and Zens, R. (2012b). Fast and scalable decoding with language model look-ahead for phrase-based statistical machine translation. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Short Papers-Volume 2*, pages 28–32. Association for Computational Linguistics.
- Yasuhara, M., Tanaka, T., Norimatsu, J.-y., and Yamamoto, M. (2013). An efficient language model using double-array structures. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 222–232, Seattle, Washington, USA. Association for Computational Linguistics.
- Zens, R. and Ney, H. (2007). Efficient phrase-table representation for machine translation with applications to online mt and speech translation. In *HLT-NAACL*, pages 492–499.
- Zens, R. and Ney, H. (2008). Improvements in dynamic programming beam search for phrase-based statistical machine translation. In *Proc. of IWSLT*.