

# Fast, Scalable Phrase-Based SMT Decoding

Anonymous ACL submission

## 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 that for the software to be optimized for their requirements, in particular, faster phrase-based decoding and more efficient utilization of modern multi-core 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, re-implementing the decoder's structure to optimize for speed, re-assessing some of the optimization assumptions made in the implementation of significant parts of the decoder, as well as exploring changes to the core search algorithm. We will compare and contrast decoder implementations, not just decoding algorithms or individual components.

We will present a phrase-based decoder that is significantly faster than the Moses baseline for single-threaded operation, and scales monotonically 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.

### 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 trade-off between speed and translation quality to the adjusted with a single parameter. Wuebker et al. (2012) 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 (Peitz et al., 2012) supports the language model look-ahead described in Wuebker et al. (2012) 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 (Moses<sup>1</sup> and Joshua<sup>2</sup>).

Jane and mtplz are single-threaded decoders, relying on external applications to parallelize operations. This has the advantage it is easier to parallelize processing across multiple servers, not just multiple cores. As we shall see, this method does not scale beyond a small number of cores when attempting to scale to multiple cores within one server.

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

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

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

The rest of the paper will be broken up into the following sections. Next, we will describe the phrase-based model and the major implementation components, with particular emphasis on decoding speed. We will then describe modifications to improve decoding speed and present results. We conclude in the last section and discuss future work.

## 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 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  $\lambda_m$  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.

### 2.1 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 have 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. Partial translations are organized into stacks where each stack holds a number of comparable hypotheses. Most phrase-based implementations place hypotheses in the same stack that have the same coverage cardinality  $|C|$ , where  $C$  is the coverage set,  $C \subseteq \{1, 2, \dots |s|\}$  of the number of source words translated.

## 2.2 Feature Functions

Features functions are the  $h_m$  in Equation 2, calculating a score for each hypothesis.

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, the first four are simple and do not leave much room for optimization. As have already been mentioned, the language model consumes a large part of decoding time but it has already been heavily optimized in previous works such as Heafield (2011) and Yasuhara et al. (2013).

Therefore, we will focus on optimizing only the phrase-table and lexicalized reordering model.

## 3 Proposed Improvements

Apart from the two feature function optimization detailed in the last section, we will investigate two other optimizations, all described below.

### 3.1 Efficient Memory Allocation

The search algorithm creates and destroy a large number of hypothesis and related objects. The memory allocation required 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. In time-critical parts of the decoder, memory will be allocated from a memory pool.

A memory pool is a large block of memory that has been pre-allocated by the operating system from which smaller blocks can be give out to the application. We will use a thread-specific memory pools to avoid locking contention during memory access. Our memory pool will be dynamic, growing in size when required but never reducing in size.

To further increase memory management speed, objects are not deleted, instead the memory pool is simply reset and reused for each input sentence. This can result in high memory usage so object queues are available for high-churn objects which allows the decoder to re-use portion of the memory pool before it is reset. The object queues are LIFO so that recently accessed memory are used, increasing CPU cache hits.

### 3.2 Stack Configurations

The most popular stack configuration for phrase-based model, as implemented in Pharaoh, Moses and Joshua, is coverage cardinality, that is, hypotheses that have translated the same number of source words are stored in the same stack. Since stack pruning occurs in each stack independently of each other, this affect the decoder's search errors.

However, there are alternatives to the coverage cardinality 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 Peitz et al. (2012) 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 config-

uration 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.

### 3.3 Phrase-Table Optimizations

For any phrase-table table of a realistic size, memory and loading time constraints requires us to use load-on-demand implementations of which Moses can make use of, each with differing performance characteristics. Figure 1 shows the time taken to decode 800,000 sentences for the fastest two implementations. From this, it appears that the Prob-

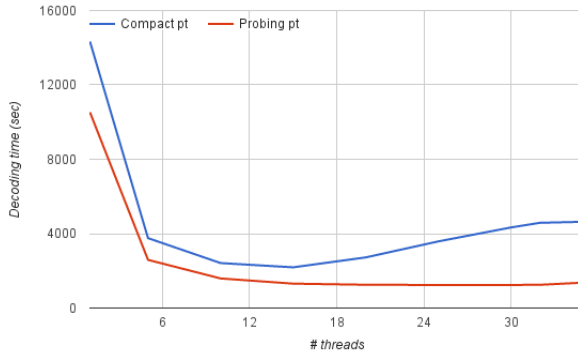


Figure 1: Moses decoding time with two different phrase-table implementations

ing phrase-table (Bogoychev, 2013) 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 periodic clearing of old rules in order to reduce memory consumption resulting in *slower* decoding, Table 1.

	No cache	Caching
Decoding time	2032	2302 (+13%)

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

We shall explore a simpler caching mechanism that creates a static cache at the start of decoding of the most likely translation rules to be used.

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 explore optimizing decoding speed by disabling compression.

### 3.4 Lexicalized Reordering Model Optimizations

Similar to the phrase-table, the lexicalized reordering model is trained on parallel data. A resultant model file is created in training which is then queried during decoding. Inevitably, the querying results in loss of 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 reordering model separately, we shall integrating it into the translation model, eliminating the need to query a separate file. This has precedent in Peitz et al. (2012) but the effect on decoding speed were not published. We will compare results with using a separate model in this paper.

## 4 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 reordering 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

	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

sentences with large models while fr-en have overly long sentences with smaller models. Where we need to compare model scores, 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 400, 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 recent<sup>3</sup> 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 1.59<sup>4</sup> and the tcmalloc library.

## 5 Results

### 5.1 Optimizing Memory

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.

<sup>3</sup>between January and May 2016

<sup>4</sup><http://boost.org/>

Hypothesis objects and the data structures it contains are re-used by maintaining a thread-specific LIFO queue of unused, recombined and pruned objects from which new hypotheses are allocated in preference to create new hypotheses in the memory pool.

Each sentence is decoded from start to finish using one worker thread. The CPU affinity of each thread is set to a specific core to minimize memory page faults as a result of threads switching cores.

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.

# threads	Moses		Our Work	
	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

Figure 2 compares the decoding time for Moses and our decoder, using the same models, parameters and test set. Our decoder is over 3 times faster with one thread, and 4.7 times faster using all cores. Like Moses, performance actually worsens after approximately 15 threads, however, the problem is not as pronounced. This gives us a better foundation on which to build further innovations for fast, multi-core decoding.

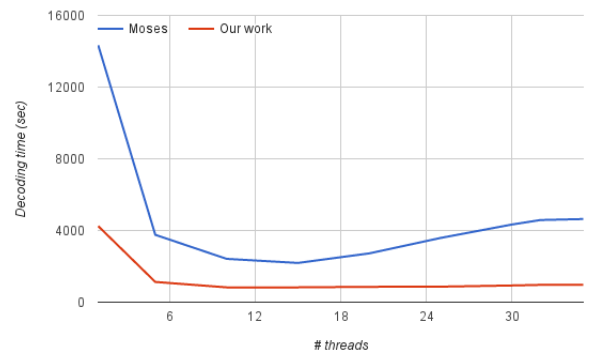


Figure 2: Decoding time of Moses and our decoder, using the same models

## 5.2 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.

```

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

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

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

```

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

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.

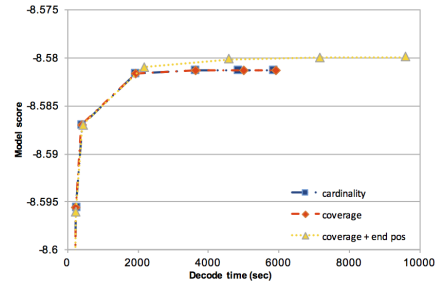


Figure 5: Trade-off between decoding time average model scores for different stack configurations

## 5.3 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 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% de-



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

crease 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

#### 5.4 Lexicalized Reordering Model

The lexicalized reordering model requires a probability distribution of the reordering 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 reordering 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 reordering 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. In fact, the decoding time with the integrated model is similar to that *without* a lexicalized reordering model<sup>5</sup>. Critically for

<sup>5</sup>Strangely, using lexicalized reordering model results in slightly faster decoding than without the model. We double checked with other datasets. We suspect this may be an arte-

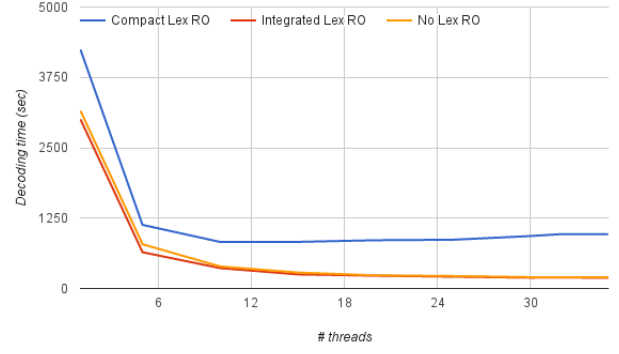


Figure 6: Decoding time with Compact Lexicalized Reordering, and integrated into a model the phrase-table

systems with multicore servers, it enables the decoder to continue to scale, making efficient use of all available cores.

#### 5.5 Scalability

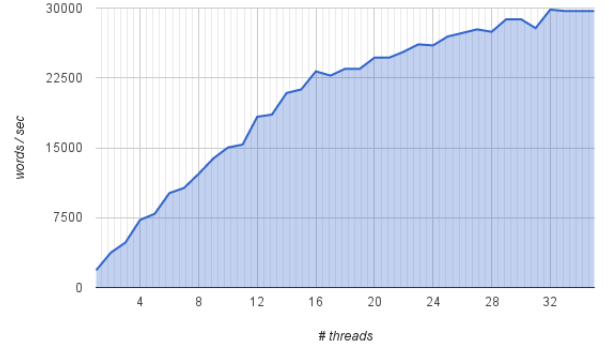


Figure 7: Our decoder’s decoding speed

Showing decoding *speed* rather than decoding time better illustrate the scalability of our approach. Figure 7 shows decoding speed against the number of threads, measured in words translated per second. 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 12.5 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 becomes slower thereafter, Figure 8.

fact of the algorithm or data structures, for example, more work is needed to resolve hypothesis recombination.

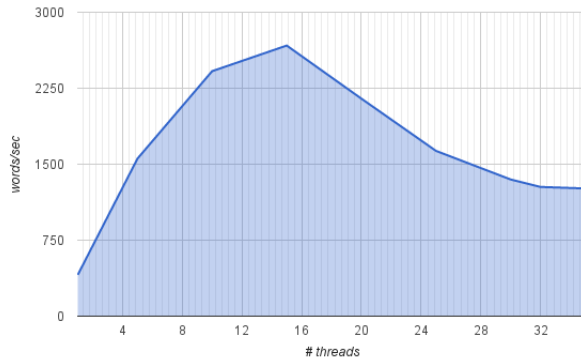


Figure 8: Moses' decoding speed

Our work is 4.5 times faster than Moses with a single-thread and 9.6 faster when all cores are used.

We mimic the use of multi-core servers by single-threaded decoders such as mtlz (Heafield et al., 2014) and Jane (Peitz et al., 2012) by running multiple processes of the Moses decoder, each running one thread. The input is distributed to the processes using a external script.

The results in ??? shows that not only does the parallelization method scales badly even compared to Moses, it is unable to run at all when more that a certain number of threads, possibly due to higher memory requirements of running multiple memory-hungry processes, or the inefficiency of distributing work to many processes. Based on this result, single-thread decoding applications are not a viable option for efficiently utilizing the multicore servers of today.

## 6 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 9.

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 10. Our decoder is 5.4 times faster than Moses with a single-thread and 14.5 faster when all cores are saturated.

## 7 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

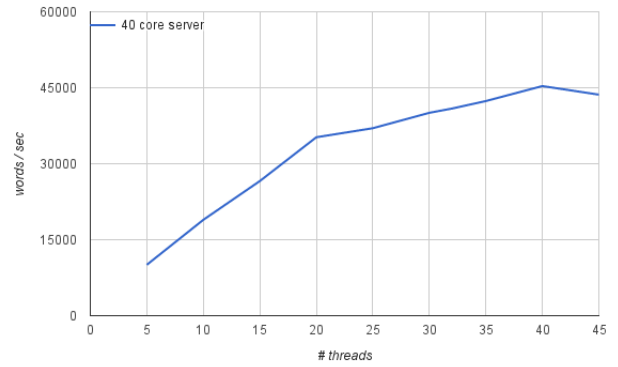


Figure 9: Decoding speed for with bigger servers

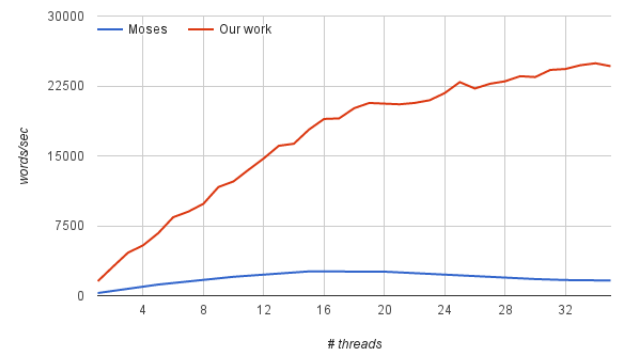


Figure 10: Decoding speed for fr-en model

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 shall also explore the underlying reasons for the scalability issues in Moses to get a better understanding where potential performance issues can arise. This has application to other algorithms beside MT decoding.

## References

- Nikolay Bogoychev. 2013. Big data: Fast, small, and read only lookup. University of Edinburgh. Undergraduate Thesis.
- Peter F. Brown, Stephen A. Della-Pietra, Vincent J. Della-Pietra, and Robert L. Mercer. 1993. The mathematics of statistical machine translation. *Computational Linguistics*, 19(2):263–313.
- David Chiang. 2007. Hierarchical phrase-based translation. *Computational Linguistics*, 33(2):201–228.



832	M. Fernández, Juan C. Pichel, José C. Cabaleiro, and	Joern Wuebker Matthias Huck Stephan Peitz, Malte	884
833	Tomás F. Pena. 2016. Boosting performance of	Nuhn, Markus Freitag Jan-Thorsten Peter Saab, and	885
834	a statistical machine translation system using dy-	Mansour Hermann Ney. 2012. June 2: Open	886
835	namic parallelism. <i>Journal of Computational Sci-</i>	source phrase-based and hierarchical statistical ma-	887
836	<i>ence</i> , 13:37–48.	chine translation. In <i>24th International Conference</i>	888
837	Sanjay Ghemawat and Paul Menage. 2009. Tcmalloc:	<i>on Computational Linguistics</i> , page 483. Citeseer.	889
838	Thread-caching malloc.		890
839		Daniel Cer Spence Green and Christopher D Manning.	891
840	Kenneth Heafield, Michael Kayser, and Christopher D.	2014. Phrasal: A toolkit for new directions in sta-	892
841	Manning. 2014. Faster Phrase-Based decoding by	tistical machine translation. In <i>Proceedings of the</i>	893
842	refining feature state. In <i>Proceedings of the Associa-</i>	<i>Ninth Workshop on Statistical Machine Translation</i> ,	894
843	<i>tion for Computational Linguistics</i> , Baltimore, MD,	pages 114–121. Citeseer.	895
844	USA, June.		896
845	Kenneth Heafield. 2011. KenLM: faster and smaller	Jörg Tiedemann. 2012. Parallel data, tools and inter-	897
846	language model queries. In <i>Proceedings of the</i>	faces in opus. In <i>LREC</i> , pages 2214–2218.	898
847	<i>EMNLP 2011 Sixth Workshop on Statistical Ma-</i>		899
848	<i>chine Translation</i> , pages 187–197, Edinburgh, Scot-	Joern Wuebker, Hermann Ney, and Richard Zens.	900
849	land, United Kingdom, July.	2012. Fast and scalable decoding with language	901
850		model look-ahead for phrase-based statistical ma-	902
851	Marcin Junczys-Dowmunt. 2012. A space-efficient	chine translation. In <i>Proceedings of the 50th An-</i>	903
852	phrase table implementation using minimal perfect	<i>nual Meeting of the Association for Computational</i>	904
853	hash functions. In Petr Sojka, Ales Hork, Ivan	<i>Linguistics: Short Papers-Volume 2</i> , pages 28–32.	905
854	Kopecek, and Karel Pala, editors, <i>15th International</i>	Association for Computational Linguistics.	906
855	<i>Conference on Text, Speech and Dialogue (TSD)</i> ,		907
856	volume 7499 of <i>Lecture Notes in Computer Science</i> ,	Makoto Yasuhara, Toru Tanaka, Jun-ya Norimatsu, and	908
857	pages 320–327. Springer.	Mikio Yamamoto. 2013. An efficient language	909
858		model using double-array structures. In <i>Proceed-</i>	910
859	Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris	<i>ings of the 2013 Conference on Empirical Methods</i>	911
860	Callison-Burch, Marcello Federico, Nicola Bertoldi,	<i>in Natural Language Processing</i> , pages 222–232,	912
861	Brooke Cowan, Wade Shen, Christine Moran,	Seattle, Washington, USA, October. Association for	913
862	Richard Zens, Chris Dyer, Ondrej Bojar, Alexan-	Computational Linguistics.	914
863	dra Constantin, and Evan Herbst. 2007. Moses:		915
864	Open source toolkit for statistical machine transla-	Richard Zens and Hermann Ney. 2007. Efficient	916
865	tion. In <i>Proceedings of the 45th Annual Meeting of</i>	phrase-table representation for machine translation	917
866	<i>the Association for Computational Linguistics Com-</i>	with applications to online mt and speech transla-	918
867	<i>panion Volume Proceedings of the Demo and Poster</i>	tion. In <i>HLT-NAACL</i> , pages 492–499.	919
868	<i>Sessions</i> , pages 177–180, Prague, Czech Republic,		920
869	June. Association for Computational Linguistics.	Richard Zens and Hermann Ney. 2008. Improvements	921
870		in dynamic programming beam search for phrase-	922
871	Zhifei Li, Chris Callison-Burch, Chris Dyer, San-	based statistical machine translation. In <i>Proc. of</i>	923
872	jeev Khudanpur, Lane Schwartz, Wren Thornton,	<i>IWSLT</i> .	924
873	Jonathan Weese, and Omar Zaidan. 2009. Joshua:		925
874	An open source toolkit for parsing-based machine		926
875	translation. In <i>Proceedings of the Fourth Workshop</i>		927
876	<i>on Statistical Machine Translation</i> , pages 135–139,		928
877	Athens, Greece, March. Association for Computa-		929
878	tional Linguistics.		930
879			931
880	Franz Josef Och, Nicola Ueffing, and Hermann Ney.		932
881	2001. An efficient A* search algorithm for statisti-		933
882	cal machine translation. In <i>Workshop on Data-</i>		934
883	<i>Driven Machine Translation at 39th Annual Meet-</i>		935
	<i>ing of the Association of Computational Linguistics</i>		
	<i>(ACL)</i> .		
	Daniel Ortiz-Martínez, Ismael García-Varea, and Fran-		
	cisco Casacuberta. 2006. Generalized stack decod-		
	ing algorithms for statistical machine translation. In		
	<i>Proceedings on the Workshop on Statistical Machine</i>		
	<i>Translation</i> , pages 64–71, New York City, June. As-		
	sociation for Computational Linguistics.		