# Beam Search Strategies for Neural Machine Translation

**Markus Freitag** and **Yaser Al-Onaizan**
IBM T.J. Watson Research Center
1101 Kitchawan Rd, Yorktown Heights, NY 10598
{freitagm,onaizan}@us.ibm.com

## Abstract

The basic concept in Neural Machine Translation (NMT) is to train a large Neural Network that maximizes the translation performance on a given parallel corpus. NMT is then using a simple left-to-right beam-search decoder to generate new translations that approximately maximize the trained conditional probability. The current beam search strategy generates the target sentence word by word from left-to-right while keeping a fixed amount of active candidates at each time step. First, this simple search is less adaptive as it also expands candidates whose scores are much worse than the current best. Secondly, it does not expand hypotheses if they are not within the best scoring candidates, even if their scores are close to the best one. The latter one can be avoided by increasing the beam size until no performance improvement can be observed. While you can reach better performance, this has the drawback of a slower decoding speed. In this paper, we concentrate on speeding up the decoder by applying a more flexible beam search strategy whose candidate size may vary at each time step depending on the candidate scores. We speed up the original decoder by up to 43% for the two language pairs German→English and Chinese→English without losing any translation quality.

## 1 Introduction

Due to the fact that Neural Machine Translation (NMT) is reaching comparable or even better performance compared to the traditional statistical machine translation (SMT) models (Jean et al., 2015; Luong et al., 2015), it has become very popular in the recent years (Kalchbrenner and Blunsom, 2013; Sutskever et al., 2014; Bahdanau et al., 2014). With the recent success of NMT, attention has shifted towards making it more practical. One of the challenges is the search strategy for extracting the best translation for a given source sentence. In NMT, new sentences are translated by a simple beam search decoder that finds a translation that approximately maximizes the conditional probability of a trained NMT model. The beam search strategy generates the translation word by word from left-to-right while keeping a fixed number (beam) of active candidates at each time step. By increasing the beam size, the translation performance can increase at the expense of significantly reducing the decoder speed. Typically, there is a saturation point at which the translation quality does not improve any more by further increasing the beam. The motivation of this work is two folded. First, we prune the search graph, thus, speed up the decoding process without losing any translation quality. Secondly, we observed that the best scoring candidates often share the same history and often come from the same partial hypothesis. We limit the amount of candidates coming from the same partial hypothesis to introduce more diversity without reducing the decoding speed by just using a higher beam.

## 2 Related Work

The original beam search for sequence to sequence models has been introduced and described by (Graves, 2012; Boulanger-Lewandowski et al., 2013) and by (Sutskever et al., 2014) for neural machine translation. (Hu et al., 2015; Mi et al., 2016) improved the beam search with a constraint softmax function which only considered a limited

word set of translation candidates to reduce the computation complexity. This has the advantage that they normalize only a small set of candidates and thus improve the decoding speed. (Wu et al., 2016) only consider tokens that have local scores that are not more than beamsize below the best token during their search. Further, the authors prune all partial hypotheses whose score are beamsize lower than the best final hypothesis (if one has already been generated). In this work, we investigate different absolute and relative pruning schemes which have successfully been applied in statistical machine translation for e.g. phrase table pruning (Zens et al., 2012).

## 3 Original Beam Search

The original beam-search strategy finds a translation that approximately maximizes the conditional probability given by a specific model. It builds the translation from left-to-right and keeps a fixed number (beam) of translation candidates with the highest log-probability at each time step. For each end-of-sequence symbol that is selected among the highest scoring candidates the beam is reduced by one and the translation is stored into a final candidate list. When the beam is zero, it stops the search and picks the translation with the highest log-probability (normalized by the number of target words) out of the final candidate list.

## 4 Search Strategies

In this section, we describe the different strategies we experimented with. In all our extensions, we first reduce the candidate list to the current beam size and apply on top of this one or several of the following pruning schemes.

**Relative Threshold Pruning.** The relative threshold pruning method discards those candidates that are far worse than the best active candidate. Given a pruning threshold $rp$ and an active candidate list $C$, a candidate $cand \in C$ is discarded if:

$$score(cand) \leq rp * \max_{c \in C}\{score(c)\} \quad (1)$$

**Absolute Threshold Pruning.** Instead of taking the relative difference of the scores into account, we just discard those candidates that are worse by a specific threshold than the best active candidate. Given a pruning threshold

$ap$ and an active candidate list $C$, a candidate $cand \in C$ is discarded if:

$$score(cand) \leq \max_{c \in C}\{score(c)\} - ap \quad (2)$$

**Relative Local Threshold Pruning.** In this pruning approach, we only consider the score $score_w$ of the last generated word and not the total score which also include the scores of the previously generated words. Given a pruning threshold $rpl$ and an active candidate list $C$, a candidate $cand \in C$ is discarded if:

$$score_w(cand) \leq rpl * \max_{c \in C}\{score_w(c)\} \quad (3)$$

**Maximum Candidates per Node** We observed that at each time step during the decoding process, most of the partial hypotheses share the same predecessor words. To introduce more diversity, we allow only a fixed number of candidates with the same history at each time step. Given a maximum candidate threshold $mc$ and an active candidate list $C$, a candidate $cand \in C$ is discarded if already $mc$ better scoring partial hyps with the same history are in the candidate list.

## 5 Experiments

For the German→English translation task, we train an NMT system based on the WMT 2016 training data (Bojar et al., 2016) (3.9M parallel sentences). For the Chinese→English experiments, we use an NMT system trained on 11 million sentences from the BOLT project.

In all our experiments, we use our in-house attention-based NMT implementation which is similar to (Bahdanau et al., 2014). For German→English, we use sub-word units extracted by byte pair encoding (Sennrich et al., 2015) instead of words which shrinks the vocabulary to 40k sub-word symbols for both source and target. For Chinese→English, we limit our vocabularies to be the top 300K most frequent words for both source and target language. Words not in these vocabularies are converted into an unknown token. During translation, we use the alignments (from the attention mechanism) to replace the unknown tokens either with potential targets (obtained from an IBM Model-1 trained on the parallel data) or with the source word itself (if no target was found) (Mi et al., 2016). We use an embedding dimension of 620 and fix the RNN GRU lay-
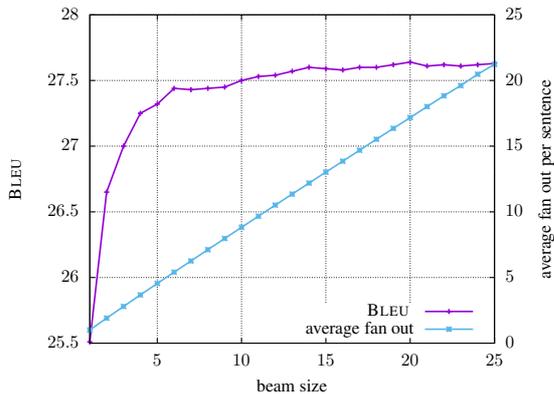
Figure 1: German→English: Original beam-search strategy with different beam sizes on newstest2014.
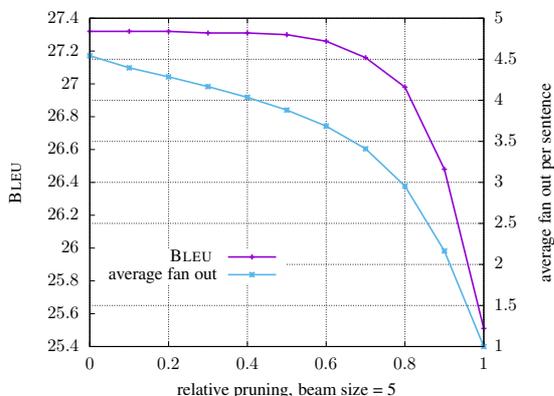


Figure 2: German→English: Different values of relative pruning measured on newstest2014.

ers to be of 1000 cells each. For the training procedure, we use SGD (Bishop, 1995) to update model parameters with a mini-batch size of 64. The training data is shuffled after each epoch.

We measure the decoding speed by two numbers. First, we compare the actual speed relative to the same setup without any pruning. Secondly, we measure the average fan out per time step. For each time step, the fan out is defined as the number of candidates we expand. Fan out has an upper bound of the size of the beam, but can be decreased either due to early stopping (we reduce the beam every time we predict a end-of-sentence symbol) or by the proposed pruning schemes. For each pruning technique, we run the experiments with different pruning thresholds and chose the largest threshold that did not degrade the translation performance based on a selection set.

In Figure 1, you can see the German→English translation performance and the average fan out per sentence for different beam sizes. Based

on this experiment, we decided to run our pruning experiments for beam size 5 and 14. The German→English results can be found in Table 1. By using the combination of all pruning techniques, we can speed up the decoding process by 13% for beam size 5 and by 43% for beam size 14 without any drop in performance. The relative pruning technique is the best working one for beam size 5 whereas the absolute pruning technique works best for a beam size 14. In Figure 2 the decoding speed with different relative pruning threshold for beam size 5 are illustrated. Setting the threshold higher than 0.6 hurts the translation performance. A nice side effect is that it has become possible to decode without any fix beam size when we apply pruning. Nevertheless, the decoding speed drops while the translation performance did not change. Further, we looked at the number of search errors introduced by our pruning schemes (number of times we prune the best scoring hypothesis). 5% of the sentences change due to search errors for beam size 5 and 9% of the sentences change for beam size 14 when using all four pruning techniques together.

The Chinese→English translation results can be found in Table 2. We can speed up the decoding process by 10% for beam size 5 and by 24% for beam size 14 without loss in translation quality. In addition, we measured the number of search errors introduced by pruning the search. Only 4% of the sentences change for beam size 5, whereas 22% of the sentences change for beam size 14.

## 6 Conclusion

The original beam search decoder used in Neural Machine Translation is very simple. It generated translations from left-to-right while looking at a fix number (beam) of candidates from the last time step only. By setting the beam size large enough, we ensure that the best translation performance can be reached with the drawback that many candidates whose scores are far away from the best are also explored. In this paper, we introduced several pruning techniques which prune candidates whose scores are far away from the best one. By applying a combination of absolute and relative pruning schemes, we speed up the decoder by up to 43% without losing any translation quality. Putting more diversity into the decoder did not improve the translation quality.

| pruning | beam size | speed up | avg fan out per sent | tot fan out per sent | newstest2014 | | newstest2015 | |
|---|---|---|---|---|---|---|---|---|
| | | | | | BLEU | TER | BLEU | TER |
| no pruning | 1 | - | 1.00 | 25 | 25.5 | 56.8 | 26.1 | 55.4 |
| no pruning | 5 | - | 4.54 | 122 | 27.3 | 54.6 | 27.4 | 53.7 |
| rp=0.6 | 5 | 6% | 3.71 | 109 | 27.3 | 54.7 | 27.3 | 53.8 |
| ap=2.5 | 5 | 5% | 4.11 | 116 | 27.3 | 54.6 | 27.4 | 53.7 |
| rpl=0.02 | 5 | 5% | 4.25 | 118 | 27.3 | 54.7 | 27.4 | 53.8 |
| mc=3 | 5 | 0% | 4.54 | 126 | 27.4 | 54.6 | 27.5 | 53.8 |
| rp=0.6,ap=2.5,rpl=0.02,mc=3 | 5 | 13% | 3.64 | 101 | 27.3 | 54.6 | 27.3 | 53.8 |
| no pruning | 14 | - | 12.19 | 363 | 27.6 | 54.3 | 27.6 | 53.5 |
| rp=0.3 | 14 | 10% | 10.38 | 315 | 27.6 | 54.3 | 27.6 | 53.4 |
| ap=2.5 | 14 | 29% | 9.49 | 279 | 27.6 | 54.3 | 27.6 | 53.5 |
| rpl=0.3 | 14 | 24% | 10.27 | 306 | 27.6 | 54.4 | 27.7 | 53.4 |
| mc=3 | 14 | 1% | 12.21 | 347 | 27.6 | 54.4 | 27.7 | 53.4 |
| rp=0.3,ap=2.5,rpl=0.3,mc=3 | 14 | 43% | 8.44 | 260 | 27.6 | 54.5 | 27.6 | 53.4 |
| rp=0.3,ap=2.5,rpl=0.3,mc=3 | - | - | 28.46 | 979 | 27.6 | 54.4 | 27.6 | 53.3 |

Table 1: Results German→English: relative pruning(rp), absolute pruning(ap), relative local pruning(rpl) and maximum candidates per node(mc). Average fan out is the average number of candidates we keep at each time step during decoding.

| pruning | beam size | speed up | avg fan out per sent | tot fan out per sent | MT08 nw | | MT08 wb | |
|---|---|---|---|---|---|---|---|---|
| | | | | | BLEU | TER | BLEU | TER |
| no pruning | 1 | - | 1.00 | 29 | 27.3 | 61.7 | 26.0 | 60.3 |
| no pruning | 5 | - | 4.36 | 137 | 34.4 | 57.3 | 30.6 | 58.2 |
| rp=0.2 | 5 | 1% | 4.32 | 134 | 34.4 | 57.3 | 30.6 | 58.2 |
| ap=5 | 5 | 4% | 4.26 | 132 | 34.3 | 57.3 | 30.6 | 58.2 |
| rpl=0.01 | 5 | 1% | 4.35 | 135 | 34.4 | 57.5 | 30.6 | 58.3 |
| mc=3 | 5 | 0% | 4.37 | 139 | 34.4 | 57.4 | 30.7 | 58.2 |
| rp=0.2,ap=5,rpl=0.01,mc=3 | 5 | 10% | 3.92 | 121 | 34.3 | 57.3 | 30.6 | 58.2 |
| no pruning | 14 | - | 11.96 | 376 | 35.3 | 57.1 | 31.2 | 57.8 |
| rp=0.2 | 14 | 3% | 11.62 | 362 | 35.2 | 57.2 | 31.2 | 57.8 |
| ap=2.5 | 14 | 14% | 10.15 | 321 | 35.2 | 56.9 | 31.1 | 57.9 |
| rpl=0.3 | 14 | 10% | 10.93 | 334 | 35.3 | 57.2 | 31.1 | 57.9 |
| mc=3 | 14 | 0% | 11.98 | 378 | 35.3 | 56.9 | 31.1 | 57.8 |
| rp=0.2,ap=2.5,rpl=0.3,mc=3 | 14 | 24% | 8.62 | 306 | 35.3 | 56.9 | 31.1 | 57.8 |
| rp=0.2,ap=2.5,rpl=0.3,mc=3 | - | - | 38.76 | 1411 | 35.2 | 57.3 | 31.1 | 57.9 |

Table 2: Results Chinese→English: relative pruning(rp), absolute pruning(ap), relative local pruning(rpl) and maximum candidates per node(mc).

# References

D. Bahdanau, K. Cho, and Y. Bengio. 2014. Neural machine translation by jointly learning to align and translate. *ArXiv e-prints* .

Christopher M Bishop. 1995. *Neural networks for pattern recognition*. Oxford university press.

Ondrej Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Varvara Logacheva, Christof Monz, et al. 2016. Findings of the 2016 conference on machine translation (wmt16). *Proceedings of WMT* .

Nicolas Boulanger-Lewandowski, Yoshua Bengio, and Pascal Vincent. 2013. Audio chord recognition with recurrent neural networks. In *ISMIR*. Citeseer, pages 335–340.

Alex Graves. 2012. Sequence transduction with recurrent neural networks. *arXiv preprint arXiv:1211.3711* .

Xiaoguang Hu, Wei Li, Xiang Lan, Hua Wu, and

Haifeng Wang. 2015. Improved beam search with constrained softmax for nmt. *Proceedings of MT Summit XV* page 297.

Sébastien Jean, Kyunghyun Cho, Roland Memisevic, and Yoshua Bengio. 2015. On using very large target vocabulary for neural machine translation. In *Proceedings of ACL*. Beijing, China, pages 1–10.

Nal Kalchbrenner and Phil Blunsom. 2013. Recurrent continuous translation models. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Seattle.

Thang Luong, Ilya Sutskever, Quoc Le, Oriol Vinyals, and Wojciech Zaremba. 2015. Addressing the rare word problem in neural machine translation. In *Proceedings of ACL*. Beijing, China, pages 11–19.

Haitao Mi, Zhiguo Wang, and Abe Ittycheriah. 2016. Vocabulary manipulation for neural machine translation. *arXiv preprint arXiv:1605.03209* .

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2015. Neural machine translation of rare words with subword units. *arXiv preprint arXiv:1508.07909* .

Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. 2014. Sequence to sequence learning with neural networks. In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*. pages 3104–3112. http://papers.nips.cc/paper/5346-sequence-to-sequence-learning-with-neural-networks.

Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. 2016. Google's neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144* .

Richard Zens, Daisy Stanton, and Peng Xu. 2012. A systematic comparison of phrase table pruning techniques. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*. Association for Computational Linguistics, pages 972–983.