

# Beam Decoding with Controlled Patience

Anonymous ACL submission

## Abstract

Text generation with beam search has proven successful in a wide range of applications. The commonly-used implementation of beam decoding follows a *first come, first served* heuristic: it keeps a set of already completed sequences over time steps and stops when the size of this set reaches the beam size. We introduce a *patience factor*, a simple modification to this decoding algorithm, that generalizes the stopping criterion and provides flexibility to the depth of search. Extensive empirical results demonstrate that the patience factor improves decoding performance of strong pre-trained models on news text summarization and machine translation over diverse language pairs, with a negligible inference slowdown. Our approach only modifies one line of code and can be thus readily incorporated in any implementation.<sup>1</sup>

## 1 Introduction

Beam search has become a dominant inference algorithm for a wide range of language generation tasks, such as machine translation (Sutskever et al., 2014; Bahdanau et al., 2015; Vaswani et al., 2017), summarization (Nallapati et al., 2016; See et al., 2017), and image captioning (Anderson et al., 2018; Li et al., 2020). Beam decoding<sup>2</sup> is an approximate, pruned version of breadth-first search that seeks the highest-probability sequence under an autoregressive (left-to-right) language generation model. In this work, we examine a popular implementation of beam decoding and propose a simple modification (one line of code) that improves the decoding performance of strong, neural language generation models (Fig. 1).

A widely-used implementation of beam language decoding (e.g., fairseq, Ott et al., 2019;

<sup>1</sup>Our codebase is available at [anonymized](#).

<sup>2</sup>In this paper, we use “beam decoding” to mean beam search applied to decoding for text generation.

---

### FCFS Beam Decoding with Controlled Patience

---

$k$ : beam size,  $M$ : maximum length,  $\mathcal{V}$ : Vocabulary  
score( $\cdot$ ): scoring function,  $p$ : patience factor.

```
1:  $B_0 \leftarrow \{(0, \text{BOS})\}$ ,  $F_0 \leftarrow \emptyset$ 
2: for  $t \in \{1, \dots, M-1\}$  :
3:    $H \leftarrow \emptyset$ ,  $F_t \leftarrow F_{t-1}$ 
4:   for  $\langle s, \mathbf{y} \rangle \in B_{t-1}$  : # Expansion.
5:     for  $y \in \mathcal{V}$  :
6:        $s \leftarrow \text{score}(\mathbf{y} \circ y)$ ,  $H.\text{add}(\langle s, \mathbf{y} \circ y \rangle)$ 
7:    $B_t \leftarrow \emptyset$ 
8:   while  $|B_t| < k$  : # Find top  $k$  w/o EOS from  $H$ .
9:      $\langle s, \mathbf{y} \rangle \leftarrow H.\text{max}()$ 
10:    if  $\mathbf{y}.\text{last}() = \text{EOS}$  :
11:       $F_t.\text{add}(\langle s, \mathbf{y} \rangle)$  # Finished hypotheses.
12:    else  $B_t.\text{add}(\langle s, \mathbf{y} \rangle)$ 
13:    if  $|F_t| \geq k \cdot p$  : # Originally,  $p=1$ .
14:      return  $F_t.\text{max}()$ 
15:     $H.\text{remove}(\langle s, \mathbf{y} \rangle)$ 
16: return  $F_t.\text{max}()$ 
```

---

Figure 1: First come, first served (FCFS) beam decoding with patience factor  $p$ . The common implementation can be considered as a special case where  $p=1$ . The highlighted line is the *only* modification that this work introduces for performance improvement.  $F_t$ : already completed sequences;  $B_t$ : beam of continuing sequences.  $H_t$ : expanded hypotheses before the top- $k$  operation. The input sequence to score is omitted.

Hugging Face’s Transformers, Wolf et al., 2020)<sup>3</sup> follows a *first come, first served* (FCFS) heuristic: when a total of  $k$  finished candidates is found ( $k$  is the beam size), it returns the best one from the  $k$  candidates and discards all of the current, unfinished  $k$  sequences in the beam. Beam size  $k$  thus determines both the breadth and depth of search. We propose a *patience factor* (Fig. 1) that decomposes these two roles and controls how many finished candidates have to be found before terminating the decoding. The patience factor generalizes the commonly-used implementation and provides flexibility in the depth of beam search by changing

<sup>3</sup>[https://github.com/pytorch/fairseq/blob/main/fairseq/sequence\\_generator.py](https://github.com/pytorch/fairseq/blob/main/fairseq/sequence_generator.py); [https://github.com/huggingface/transformers/blob/master/src/transformers/generation\\_utils.py](https://github.com/huggingface/transformers/blob/master/src/transformers/generation_utils.py).

the stopping criterion.

We apply the one-line modification to strong off-the-shelf transformer models without any change to the trained models for machine translation (Tang et al., 2021) and text summarization (Lewis et al., 2020). Our experiments demonstrate that our method outperforms the original algorithm on the CNN/Dailymail (Hermann et al., 2015) and XSUM (Narayan et al., 2018) news summarization tasks and the WMT 2020/2021 machine translation tasks (Barrault et al., 2020; Akhbardeh et al., 2021) across diverse language pairs. Further, the introduction of the patience factor only results in a negligible inference slowdown, confirming its practical advantage in downstream applications.

Our analysis shows that, while the performance gain is sensitive to hyperparameters of beam decoding (beam size and length penalty; Johnson et al., 2017), the patience factor is consistently beneficial. Moreover, we extensively compare our results with the *vanilla* implementation of beam search that much prior work assumes (Meister et al., 2020b; Stahlberg and Byrne, 2019, *inter alia*). Empirically, we found that the vanilla algorithm performs competitively with FCFS on machine translation but substantially underperforms on summarization. The FCFS beam decoding with our patience factor is thus a simple yet effective algorithm for both language generation tasks.

## 2 Beam Decoding with Patience

**Vanilla and FCFS Implementations** Beam decoding has been applied to sequence-to-sequence models (Graves, 2012; Boulanger-Lewandowski et al., 2013a,b), and it is now used in many state-of-the-art systems for language generation tasks (Zhang et al., 2020, 2021; Tran et al., 2021; Raffel et al., 2020, *inter alia*). Figs. 1 and 2 describe its two major implementations. They differ primarily in the treatment of finished sequences with the EOS symbol at the end: FCFS collects finished sequences in a *first come, first served* manner and removes them from the beam (Line 11, Fig. 1), whereas the vanilla version finds the top  $k$  sequences, including both finished and unfinished sequences (Line 5 in Fig. 2). While often unspecified in the literature, our later experiments in §3.2 will show that this difference can affect the downstream performance substantially, especially on news text

<sup>4</sup>[https://www.tensorflow.org/addons/api\\_docs/python/tfa/seq2seq/BeamSearchDecoder](https://www.tensorflow.org/addons/api_docs/python/tfa/seq2seq/BeamSearchDecoder).

### Vanilla Beam Decoding

---

$k$ : beam size,  $M$ : maximum length,  
 $\mathcal{V}$ : Vocabulary,  $\text{score}(\cdot)$ : scoring function.

```

1:  $B_0 \leftarrow \{ \langle 0, \text{BOS} \rangle \}$ 
2: for  $t \in \{1, \dots, M-1\}$  :
3:   for  $\langle s, \mathbf{y} \rangle \in B_{t-1}$  :
4:     if  $\mathbf{y}.\text{last}() = \text{EOS}$  :
5:        $H.\text{add}(\langle s, \mathbf{y} \rangle)$ 
6:     continue
7:     for  $y \in \mathcal{V}$  :
8:        $s \leftarrow \text{score}(\mathbf{y} \circ y)$ ,  $H.\text{add}(\langle s, \mathbf{y} \circ y \rangle)$ 
9:    $B_t \leftarrow \emptyset$ 
10:  while  $|B_t| < k$  : # Find top  $k$  from  $H$ .
11:     $\langle s, \mathbf{y} \rangle \leftarrow H.\text{max}()$ ,  $B_t.\text{add}(\langle s, \mathbf{y} \rangle)$ 
12:     $H.\text{remove}(\langle s, \mathbf{y} \rangle)$ 
13:  if  $\mathbf{y}.\text{last}() = \text{EOS}, \forall \mathbf{y} \in B_t$  : # All finished.
14:    return  $B_t.\text{max}()$ 
15: return  $B_t.\text{max}()$ 

```

---

Figure 2: The vanilla version of beam decoding. The top- $k$  operation is applied over  $H$ , the union of the finished and continuing sequences. This is implemented, for example, in the TensorFlow Addons library (Abadi et al., 2015).<sup>4</sup> See also Stahlberg and Byrne (2019); Meister et al. (2020b).

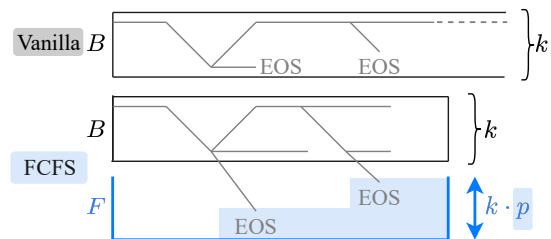


Figure 3: FCFS with patience factor  $p$  vs. vanilla beam decoding.  $k$  denotes the beam size. FCFS stores finished sentences in  $F$ , but they stay in (and later may fall off from) beam  $B$  during vanilla decoding.  $k \cdot p$  determines the size of  $F$ . The illustration of beam decoding here is inspired by Huang et al. (2012).

summarization.

Further comparing Figs. 1 and 2, we see their difference in terms of the breadth and depth of search. Given the same beam size  $k$ , FCFS has a wider breadth since it collects  $k$  unfinished sequences at every step regardless of how many sequences are finished with the EOS symbol.<sup>5</sup> The vanilla algorithm decodes until all top- $k$  sequences are finished (Line 13, Fig. 2), and therefore it tends to result in deeper search. FCFS, in contrast, terminates when a total of  $k$  finished sequences is found.

**Patience Factor for FCFS** Beam size  $k$  in FCFS

<sup>5</sup>In practice, this is implemented by taking the top  $2k$  sequences at every step. We find at most  $k$  EOS symbols, so there are always at least  $k$  unfinished sequences. See <https://github.com/huggingface/transformers/>.

thus controls both the breadth and stopping criterion (i.e., depth) of search. We introduce the patience factor (Line 13, Fig. 1) that relaxes this assumption and separates the stopping criterion from the search breadth. Fig. 3 illustrates this patience factor as well as the difference between the FCFS and vanilla algorithms. The one-line change generalizes FCFS ( $p=1$ ) and adds flexibility. We will show that this flexibility is beneficial on machine translation and summarization (§3.2).

### 3 Experiments

We present extensive comparisons of beam decoding variants on text summarization and machine translation over a wide range of language pairs. Our simple addition of the patience factor improves performance across the board.

#### 3.1 Experimental Setup

We evaluate four decoding algorithms on machine translation and summarization: **greedy**, **vanilla**, **FCFS**, and **FCFS with the patience factor**. For machine translation, we use multilingual BART (Tang et al., 2021), a strong, pretrained transformer model,<sup>6</sup> and WMT 2020/2021 news test data (Barrault et al., 2020; Akhbardeh et al., 2021) for four diverse language pairs (eight directions): WMT 2020 for EN↔PL (Polish) and 2021 for EN↔DE (German), EN↔JA (Japanese), and EN↔ZH (Chinese). We apply beam decoding with the same hyperparameters as Tang et al. (2021): beam size 5 and length penalty 1. We measure performance with the COMET score (Rei et al., 2020a,b), a state-of-the-art evaluation metric based on multilingual contextual representations. For summarization, we experiment with the CNN/Dailymail (CNNDM, Hermann et al., 2015) and XSUM (Narayan et al., 2018) datasets. We apply the off-the-shelf BART models (Lewis et al., 2020) that are fine-tuned on each dataset.<sup>7</sup> Performance is measured with ROUGE scores (Lin, 2004). We follow the original setting in Lewis et al. (2020): beam sizes 4 and 6 and length penalty 2 and 1 for CNNDM and XSUM, respectively. More experimental details are described in Appendix §A.

We experiment with the same patience factor on all datasets for each task, based on our preliminary development:  $p=2$  for machine translation and

$p=0.5$  for summarization. Here we avoid additional effort and demonstrate the practical value of our simple modification. We present detailed sensitivity analysis over  $p$  in §3.3.

#### 3.2 Results

Seen in Table 1 are results from our experiments. FCFS with the patience factor outperforms the widely-used FCFS algorithm across the board; e.g., 53.0 vs. 52.1 on EN→PL. Particularly noteworthy are the performance gains on the two summarization datasets; e.g., 31.2 vs. 30.3 ROUGE-L on CNNDM. Comparing vanilla decoding and FCFS, we see that the former outperforms the latter (and is competitive with or slightly better than FCFS w/  $p$ ) on machine translation but underperforms substantially on summarization; e.g., 34.4 vs. 33.1 ROUGE-L on XSUM. Vanilla decoding even performs worse than greedy decoding in many cases. We suspect this performance degradation on summarization might be a reason why FCFS is used instead of vanilla decoding in popular libraries.

#### 3.3 Analysis

Here we use the standard dev. split from the XSUM dataset and news test 2020 EN→DE and ZH→EN data. We fixed the value of  $p$  for each task so far, but Fig. 4 explores varying patience factors and their effects on the performance (A: EN→DE; B: XSUM) and the inference speed (C). The translation performance improves with larger patience factors with diminishing gains. On the other hand, summarization benefits more from patience factors smaller than the original value of 1, possibly due to issues in the scoring function (Wiseman and Rush, 2016) or ROUGE evaluations (Nenkova, 2006) and the nature of the summarization task that aims to generate concise text. Note, however, that we see consistent patterns with ROUGE from COMET (Rei et al., 2020b), which achieves the highest correlation to human judgment on CNNDM (Kasai et al., 2022a; see Table 3 in the appendix). Regardless, our patience factor provides useful flexibility for any generation task.

As expected, generation slows down as  $p$  increases (Fig. 4C). The inference slowdown from around  $p=2$  is still negligible, again showing the practicality of our method. Fig. 5 explores the performance gains from the patience factor over varying beam sizes. The amount of improvement changes, but the patience factor is generally beneficial. We see similar trends for various values of

<sup>6</sup><https://github.com/pytorch/fairseq/tree/main/examples/multilingual#mbart50-models>.

<sup>7</sup><https://github.com/pytorch/fairseq/tree/main/examples/bart>.

Algorithm	WMT 2020/2021 Machine Translation ( $p=2$ )								Summarization ( $p=0.5$ )					
	EN↔DE		EN↔JA		EN↔PL		EN↔ZH		CNNDM			XSUM		
	→	←	→	←	→	←	→	←	R-2	R-3	R-L	R-2	R-3	R-L
Greedy	43.7	66.2	33.6	9.5	46.0	53.5	32.5	23.5	21.1	11.9	30.7	19.8	10.7	34.3
Vanilla	48.2	66.3	<b>38.7</b>	<b>15.7</b>	52.7	58.2	<b>33.9</b>	29.9	19.2	11.0	28.0	19.5	10.7	33.1
FCFS	47.9	66.2	38.0	15.0	52.1	58.1	33.7	29.6	20.4	11.6	30.3	20.4	11.4	34.4
FCFS w/ $p$	<b>48.3</b>	<b>66.4</b>	38.4	15.6	<b>53.0</b>	<b>58.4</b>	33.8	<b>30.2</b>	<b>21.4</b>	<b>12.4</b>	<b>31.2</b>	<b>21.0</b>	<b>11.8</b>	<b>35.4</b>

Table 1: We evaluate the four inference algorithms on the machine translation and news summarization test data with the COMET score (Rei et al., 2020b) and ROUGE scores (ROUGE-2/3/L), respectively. FCFS w/  $p$  indicates our FCSF algorithm with the patience factor ( $p=2$  for machine translation and  $p=0.5$  for summarization). COMET uses crosslingual contextual representations from XLM-RoBERTa (Conneau et al., 2020) and has shown to have significantly higher correlation with expert human judgment than alternatives (Mathur et al., 2020b; Kasai et al., 2022a) like BLEU (Papineni et al., 2002). Nonetheless, we see consistent patterns from BLEU (Appendix §B). For CNNDM, we used 100 test articles with 10 human-written references each from Kryscinski et al. (2019).

the length penalty (see Fig. 6 in the appendix).

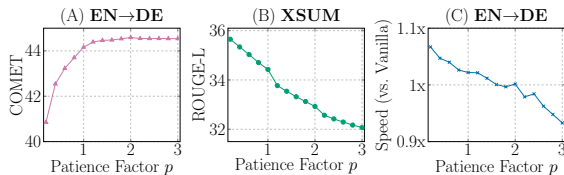


Figure 4: Effects of varying patience factors  $p$  on the dev. score (A and B) and inference speed (C). The inference speed is measured with batch size 20, relative to the vanilla decoding algorithm on the same single Nvidia A100-SXM GPU. Other languages pairs were similar to EN→DE (A). CNNDM also had similar trends to XSUM (B).

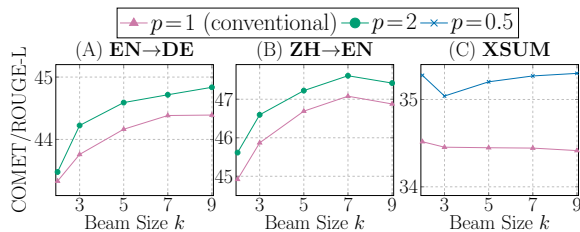


Figure 5: Effects of controlled patience on the dev. data over varying beam sizes. The length penalty value is 1. We evaluate with COMET for machine translation and ROUGE-L for XSUM summarization.

## 4 Further Related Work

**Stopping Criteria for Beam Decoding** The patience factor changes the stopping criterion and adds flexibility in the search depth of the common beam search algorithm. Similarly, several prior works studied stopping criteria to improve machine translation (Huang et al., 2017; Yang et al., 2018; Ma et al., 2019). Our machine translation experiments are consistent with their findings: stopping

criteria that yield accurate search improve performance. In the case of summarization, however, we observed that less patient and thus *less accurate* search can improve ROUGE scores.

**Breadth of Beam Decoding** Much prior work explored downstream effects of the search breadth (Koehn and Knowles, 2017; Murray and Chiang, 2018; Ott et al., 2018; Cohen and Beck, 2019; Stahlberg and Byrne, 2019, *inter alia*). Beam decoding with larger beam sizes can find sequences with higher scores but lead to performance degradation (often called the *beam search curse*; Yang et al., 2018). Recent work (Meister et al., 2020a) argued that beam decoding with small beams introduces bias that is related to the uniform information density of human-produced text (Levy, 2005). Freitag and Al-Onaizan (2017) proposed a method to adaptively shrink the beam width based on the partial scores to speed up inference. This work focused on the stopping criteria (i.e., depth) and separated them from the breadth of the commonly-used beam decoding.

## 5 Conclusion

We introduced the patience factor that generalizes the widespread implementation of beam text decoding. Our extensive experiments showed that the patience factor improves the generation performance of strong, off-the-shelf models on machine translation and summarization with an insignificant slowdown in generation. As it only requires a minimal change in code, we hope that many researchers and practitioners of language generation will benefit from our simple yet effective modification.

## 250 Limitations and Ethical Considerations

251 We evaluated our decoding method both on ma-  
252 chine translation and news summarization. Our  
253 machine translation experiments span diverse lan-  
254 guages, including morphologically rich languages  
255 (e.g., Japanese and Polish) and languages with non-  
256 Latin scripts (e.g., Japanese and Chinese). Nonethe-  
257 less, our summarization experiments are limited to  
258 English and the news domain mainly due to our  
259 budget constraints. There are also many other lan-  
260 guage generation tasks for which our method can  
261 be useful. Since our improvement only requires  
262 one line of code, we hope that practitioners will  
263 implement it for the domain and the task of their  
264 interest and further assess how our decoding algo-  
265 rithm performs over a wider range of applications.

266 Evaluating language generation remains a chal-  
267 lenging research problem. We carefully set up our  
268 experiments to mitigate potential evaluation issues.  
269 The WMT 2020/2021 test data consist only of news  
270 text written in the original language, in contrast to  
271 the test data from WMT 2018 (Bojar et al., 2018)  
272 or earlier. For example, the WMT 2021 EN→DE  
273 and DE→EN test data come from completely dif-  
274 ferent documents. This avoids the *translationese*  
275 *effect* that would overestimate the translation per-  
276 formance due to the simplicity of translated text  
277 (Graham et al., 2020). Moreover, some language  
278 pairs in the WMT 2020 and 2021 test data have mul-  
279 tiple references per instance, which increases the  
280 correlation of automatic evaluations with human  
281 judgment (Kasai et al., 2022a). We presented re-  
282 sults using automatic metrics from recent work (Rei  
283 et al., 2020b) as well as conventional, n-gram over-  
284 lap metrics (Papineni et al., 2002; Lin, 2004). Re-  
285 cent automatic metrics have shown to have higher  
286 correlation with human judgements, but human  
287 judgments are sometimes inconsistent, especially  
288 when crowdsourced (Clark et al., 2021; Kasai et al.,  
289 2022b). Since our decoding method is a generaliza-  
290 tion of the widely-used beam search algorithm, we  
291 hope that it will be tested and used in real-world  
292 systems of language generation.

## 293 References

294 Martín Abadi, Ashish Agarwal, Paul Barham, Eugene  
295 Brevdo, Zhifeng Chen, Craig Citro, Greg S. Cor-  
296 rado, Andy Davis, Jeffrey Dean, Matthieu Devin,  
297 Sanjay Ghemawat, Ian Goodfellow, Andrew Harp,  
298 Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal  
299 Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh

Levenberg, Dan Mané, Rajat Monga, Sherry Moore,  
Derek Murray, Chris Olah, Mike Schuster, Jonathon  
Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar,  
Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan,  
Fernanda Viégas, Oriol Vinyals, Pete Warden, Mar-  
tin Wattenberg, Martin Wicke, Yuan Yu, and Xiao-  
qiang Zheng. 2015. [TensorFlow: Large-scale ma-  
chine learning on heterogeneous systems](#). Software  
available from tensorflow.org. 308

Farhad Akhbardeh, Arkady Arkhangorodsky, Mag-  
dalena Biesialska, Ondřej Bojar, Rajen Chatter-  
jee, Vishrav Chaudhary, Marta R. Costa-jussa,  
Cristina España-Bonet, Angela Fan, Christian Fed-  
ermann, Markus Freitag, Yvette Graham, Ro-  
man Grundkiewicz, Barry Haddow, Leonie Harter,  
Kenneth Heafield, Christopher Homan, Matthias  
Huck, Kwabena Amponsah-Kaakyire, Jungo Kasai,  
Daniel Khashabi, Kevin Knight, Tom Kocmi, Philipp  
Koehn, Nicholas Lourie, Christof Monz, Makoto  
Morishita, Masaaki Nagata, Ajay Nagesh, Toshiaki  
Nakazawa, Matteo Negri, Santanu Pal, Allahsera Au-  
guste Tapo, Marco Turchi, Valentin Vydrin, and Mar-  
cos Zampieri. 2021. [Findings of the 2021 conference  
on machine translation \(WMT21\)](#). In *Proc. of WMT*. 323

Peter Anderson, Xiaodong He, Chris Buehler, Damien  
Teney, Mark Johnson, Stephen Gould, and Lei Zhang.  
2018. [Bottom-up and top-down attention for image  
captioning and visual question answering](#). In *Proc.  
of CVPR*. 328

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Ben-  
gio. 2015. [Neural machine translation by jointly  
learning to align and translate](#). In *Proc. of ICLR*. 331

Loïc Barrault, Magdalena Biesialska, Ondřej Bo-  
jar, Marta R. Costa-jussa, Christian Federmann,  
Yvette Graham, Roman Grundkiewicz, Barry Had-  
dow, Matthias Huck, Eric Joanis, Tom Kocmi,  
Philipp Koehn, Chi-kiu Lo, Nikola Ljubešić, Christof  
Monz, Makoto Morishita, Masaaki Nagata, Toshi-  
aki Nakazawa, Santanu Pal, Matt Post, and Marcos  
Zampieri. 2020. [Findings of the 2020 conference on  
machine translation \(WMT20\)](#). In *Proc. of WMT*. 340

Ondřej Bojar, Christian Federmann, Mark Fishel, Yvette  
Graham, Barry Haddow, Philipp Koehn, and Christof  
Monz. 2018. [Findings of the 2018 conference on  
machine translation \(WMT18\)](#). In *Proc. of WMT*. 344

Nicolas Boulanger-Lewandowski, Yoshua Bengio, and  
Pascal Vincent. 2013a. [Audio chord recognition with  
recurrent neural networks](#). In *Proc. of ISMIR*. 347

Nicolas Boulanger-Lewandowski, Yoshua Bengio, and  
Pascal Vincent. 2013b. [High-dimensional sequence  
transduction](#). In *Proc. of ICASSP*. 350

Elizabeth Clark, Tal August, Sofia Serrano, Nikita  
Haduong, Suchin Gururangan, and Noah A. Smith.  
2021. [All that’s ‘human’ is not gold: Evaluating  
human evaluation of generated text](#). In *Proc. of ACL*. 354

355	Eldan Cohen and Christopher Beck. 2019. <a href="#">Empirical analysis of beam search performance degradation in neural sequence models</a> . In <i>Proc. of ICML</i> .	409
356		410
357		411
358	Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. <a href="#">Unsupervised cross-lingual representation learning at scale</a> . In <i>Proc. of ACL</i> .	412
359		413
360		414
361		415
362		416
363		417
364	Sergey Edunov, Myle Ott, Marc’Aurelio Ranzato, and Michael Auli. 2020. <a href="#">On the evaluation of machine translation systems trained with back-translation</a> . In <i>Proc. of ACL</i> .	418
365		419
366		420
367		421
368	Markus Freitag and Yaser Al-Onaizan. 2017. <a href="#">Beam search strategies for neural machine translation</a> . In <i>Proc. of NGT</i> .	422
369		423
370		
371	Yvette Graham, Barry Haddow, and Philipp Koehn. 2020. <a href="#">Translationese in machine translation evaluation</a> .	424
372		425
373		426
374	Alex Graves. 2012. <a href="#">Sequence transduction with recurrent neural networks</a> . In <i>Representation Learning Workshop</i> .	427
375		428
376		429
377	Karl Moritz Hermann, Tomáš Kociský, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. <a href="#">Teaching machines to read and comprehend</a> . In <i>Proc. of NeurIPS</i> .	430
378		431
379		432
380		433
381	Liang Huang, Suphan Fayong, and Yang Guo. 2012. <a href="#">Structured perceptron with inexact search</a> . In <i>Proc. of NAACL</i> .	434
382		435
383		436
384	Liang Huang, Kai Zhao, and Mingbo Ma. 2017. <a href="#">When to finish? optimal beam search for neural text generation (modulo beam size)</a> . In <i>Proc. of EMNLP</i> .	437
385		438
386		439
387	Melvin Johnson, Mike Schuster, Quoc V. Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat, Fernanda Viégas, Martin Wattenberg, Greg Corrado, Macduff Hughes, and Jeffrey Dean. 2017. <a href="#">Google’s multilingual neural machine translation system: Enabling zero-shot translation</a> . <i>TACL</i> .	440
388		441
389		
390		
391		
392		
393	Jungo Kasai, Keisuke Sakaguchi, Ronan Le Bras, Lavinia Dunagan, Jacob Morrison, Alexander R. Fabri, Yejin Choi, and Noah A. Smith. 2022a. <a href="#">Bidimensional leaderboards: Generate and evaluate language hand in hand</a> . In <i>Proc. of NAACL</i> .	442
394		443
395		444
396		445
397		446
398	Jungo Kasai, Keisuke Sakaguchi, Lavinia Dunagan, Jacob Morrison, Ronan Le Bras, Yejin Choi, and Noah A. Smith. 2022b. <a href="#">Transparent human evaluation for image captioning</a> . In <i>Proc. of NAACL</i> .	447
399		448
400		
401		
402	Philipp Koehn and Rebecca Knowles. 2017. <a href="#">Six challenges for neural machine translation</a> . In <i>Proc. of NGT</i> .	449
403		450
404		451
405	Wojciech Kryscinski, Nitish Shirish Keskar, Bryan McCann, Caiming Xiong, and Richard Socher. 2019. <a href="#">Neural text summarization: A critical evaluation</a> . In <i>Proc. of EMNLP</i> .	452
406		453
407		454
408		455
	Roger Levy. 2005. <i>Probabilistic Models of Word Order and Syntactic Discontinuity</i> . Ph.D. thesis, Stanford University.	456
		457
		458
	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. <a href="#">BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension</a> . In <i>Proc. of ACL</i> .	459
		460
		461
		462
		463
		464
		465
		466
		467
		468
		469
		470
		471
		472
		473
		474
		475
		476
		477
		478
		479
		480
		481
		482
		483
		484
		485
		486
		487
		488
		489
		490
		491
		492
		493
		494
		495
		496
		497
		498
		499
		500

459	Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan,	Yilin Yang, Liang Huang, and Mingbo Ma. 2018. <a href="#">Break-</a>	513
460	Sam Gross, Nathan Ng, David Grangier, and Michael	ing the beam search curse: A study of (re-)scoring	514
461	Auli. 2019. <a href="#">fairseq: A fast, extensible toolkit for</a>	methods and stopping criteria for neural machine	515
462	sequence modeling. In <i>Proc. of NAACL: Demonstra-</i>	translation. In <i>Proc. of EMNLP</i> .	516
463	<i>tions</i> .		
464	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Pe-	517
465	Jing Zhu. 2002. <a href="#">BLEU: a method for automatic eval-</a>	ter J. Liu. 2020. <a href="#">PEGASUS: Pre-training with ex-</a>	518
466	uation of machine translation. In <i>Proc. of ACL</i> .	tracted gap-sentences for abstractive summarization.	519
467		In <i>Proc. of ICML</i> .	520
468	Matt Post. 2018. <a href="#">A call for clarity in reporting BLEU</a>	Pengchuan Zhang, Xiujun Li, Xiaowei Hu, Jianwei	521
469	scores. In <i>Proc. of WMT</i> .	Yang, Lei Zhang, Lijuan Wang, Yejin Choi, and Jian-	522
470	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	feng Gao. 2021. <a href="#">VinVL: Making visual representa-</a>	523
471	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,	tions matter in vision-language models. In <i>Proc. of</i>	524
472	Wei Li, and Peter J. Liu. 2020. <a href="#">Exploring the limits</a>	<i>CVPR</i> .	525
473	<a href="#">of transfer learning with a unified text-to-text trans-</a>		
474	former. <i>JMLR</i> .		
475	Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon		
476	Lavie. 2020a. <a href="#">COMET: A neural framework for MT</a>		
477	evaluation. In <i>Proc. of EMNLP</i> .		
478	Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon		
479	Lavie. 2020b. <a href="#">Unbabel’s participation in the WMT20</a>		
480	<a href="#">metrics shared task</a> . In <i>Proc. of WMT</i> .		
481	Abigail See, Peter J. Liu, and Christopher D. Manning.		
482	2017. <a href="#">Get to the point: Summarization with pointer-</a>		
483	<a href="#">generator networks</a> . In <i>Proc. of ACL</i> .		
484	Felix Stahlberg and Bill Byrne. 2019. <a href="#">On NMT search</a>		
485	<a href="#">errors and model errors: Cat got your tongue?</a> In		
486	<i>Proc. of EMNLP</i> .		
487	Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014.		
488	<a href="#">Sequence to sequence learning with neural networks</a> .		
489	In <i>Proc. of NeurIPS</i> .		
490	Yuqing Tang, Chau Tran, Xian Li, Peng-Jen Chen, Na-		
491	man Goyal, Vishrav Chaudhary, Jiatao Gu, and An-		
492	gela Fan. 2021. <a href="#">Multilingual translation from denois-</a>		
493	<a href="#">ing pre-training</a> . In <i>Findings of ACL</i> .		
494	Chau Tran, Shruti Bhosale, James Cross, Philipp Koehn,		
495	Sergey Edunov, and Angela Fan. 2021. <a href="#">Facebook</a>		
496	<a href="#">AI’s WMT21 news translation task submission</a> . In		
497	<i>Proc. of WMT</i> .		
498	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob		
499	Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz		
500	Kaiser, and Illia Polosukhin. 2017. <a href="#">Attention is all</a>		
501	<a href="#">you need</a> . In <i>Proc. of NeurIPS</i> .		
502	Sam Wiseman and Alexander M. Rush. 2016. <a href="#">Sequence-</a>		
503	<a href="#">to-sequence learning as beam-search optimization</a> . In		
504	<i>Proc. of EMNLP</i> .		
505	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien		
506	Chaumond, Clement Delangue, Anthony Moi, Pier-		
507	ric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz,		
508	Joe Davison, Sam Shleifer, Patrick von Platen, Clara		
509	Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven		
510	Le Scao, Sylvain Gugger, Mariama Drame, Quentin		
511	Lhoest, and Alexander Rush. 2020. <a href="#">Transformers:</a>		
512	<a href="#">State-of-the-art natural language processing</a> . In <i>Proc.</i>		
	<i>of EMNLP: System Demonstrations</i> .		

# Appendices

Hyperparameter	Value
<b>WMT Machine Translation (All Pairs)</b>	
beam size	5
length penalty	1
<b>CNNNDM Summarization</b>	
beam size	4
length penalty	2
max-len-b	140
min-len	55
no-repeat-ngram-size	3
<b>XSUM Summarization</b>	
beam size	6
length penalty	1
max-len-b	60
min-len	10
no-repeat-ngram-size	3

Table 2: Beam decoding hyperparameters. We generally followed prior work: Tang et al. (2021) for machine translation and Lewis et al. (2020) for CNNNDM and XSUM summarization.

## A Beam Decoding Hyperparameters

Table 2 shows the beam decoding hyperparameters in our experiments. We generally follow the original settings of the pretrained, off-the-shelf models (Tang et al., 2021; Lewis et al., 2020).

## B Additional Results

Table 3 reports BLEU (Papineni et al., 2002) and COMET (Rei et al., 2020b) scores for the machine translation and summarization experiments, respectively. We use the sacreBLEU implementation for BLEU (Post, 2018). Note that much recent work (Mathur et al., 2020a; Kasai et al., 2022a,b; Edunov et al., 2020, *inter alia*) found poor correlation between BLEU scores and human judgment for evaluating strong language generation models. COMET is an automatic metric for machine translation that uses crosslingual contextual representations from XLM-RoBERTa (Conneau et al., 2020), but it can be used *monolingually* for evaluating summarization as well (Kasai et al., 2022a).

Fig. 6 explores the performance gains from the patience factor over varying length penalty values. Consistent with the trends from various beam sizes (Fig. 5), the amount of improvement changes, but the patience factor is generally beneficial.

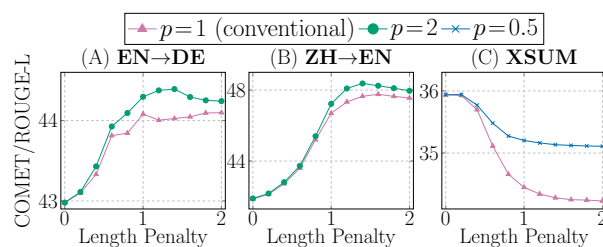


Figure 6: Effects of controlled patience on the dev. data over varying length penalty values. The beam sizes are all 5. We evaluate with COMET for machine translation and ROUGE-L for XSUM summarization.



Algorithm	WMT 2020 and 2021 Machine Translation (BLEU)								Summarization	
	EN↔DE		EN↔JA		EN↔PL		EN↔ZH		CNNDM	XSUM
	→	←	→	←	→	←	→	←	COMET	COMET
Greedy	42.9	46.6	20.2	17.4	19.8	30.7	31.2	21.7	<b>1.6</b>	0.1
Vanilla	<b>45.1</b>	48.4	21.6	19.7	<b>21.1</b>	<b>32.5</b>	<b>32.5</b>	23.6	-5.5	-1.6
FCFS	45.0	48.4	21.3	19.5	21.0	32.4	<b>32.6</b>	23.4	-4.2	2.2
FCFS w/ $p$	45.0	<b>48.5</b>	<b>21.7</b>	<b>19.8</b>	<b>21.1</b>	<b>32.5</b>	32.3	<b>23.7</b>	-1.1	<b>2.5</b>

Table 3: We evaluate the four decoding algorithms on machine translation and summarization and report BLEU (Papineni et al., 2002) and COMET (Rei et al., 2020b) scores here. FCFS w/  $p$  indicates our FCSF algorithm with the patience factor ( $p=2$  for machine translation and  $p=0.5$  for summarization). COMET is an automatic metric for machine translation that uses crosslingual contextual representations from XLM-RoBERTa (Conneau et al., 2020), but it can be used for evaluating summarization as well (Kasai et al., 2022a).