

# Self-training Reduces Flicker in Retranslation-based Simultaneous Translation

Anonymous ACL submission

## Abstract

In simultaneous translation, the *retranslation* approach has the advantage of requiring no modifications to the inference engine. However, in order to reduce the undesirable flicker in the output, previous work has resorted to increasing the latency through masking, and introducing specialised inference, thus losing the simplicity of the approach. In this work, we will show that self-training improves the flicker-latency tradeoff, whilst maintaining similar translation quality to the original. Our analysis indicates that self-training reduces flicker by controlling monotonicity. Self-training can be combined with biased beam search to further improve the flicker-latency tradeoff.

## 1 Introduction

Simultaneous machine translation systems, which process their input word by word instead of sentence by sentence, must strike a balance between producing output immediately (and so reducing quality because of incomplete input) and waiting for further input (and so increasing latency). A good simultaneous translation system will provide a pareto-optimal tradeoff between quality and latency. A straightforward way of doing simultaneous translation is *retranslation* (Niehues et al., 2016), which has the advantage that it can be used with an unmodified machine translation (MT) inference engine, and can perform better than the alternative, streaming-based approaches (Arivazhagan et al., 2020b). The disadvantage is that retranslation may change previous output causing *flicker*, leading to a poor user experience, and so flicker needs to be balanced with latency and quality.

We argue that flickering is caused by two different (but related) issues: (i) lexical instability of the translation – the system “changes its mind” as more source is revealed, swapping one word

for another<sup>1</sup>; (ii) non-monotonicity of the translation – the system favours a non-monotonic translation, which means it needs high latency in order to avoid flicker. Some of this instability and non-monotonicity is necessary – forced by syntactic differences between source and target, and lack of information in the prefixes – but some is due to arbitrary choices of the model so we aim to reduce these as far as possible.

In non-autoregressive translation (NAT), a related problem, known as the “multimodality” problem (Gu et al., 2018), has been addressed using knowledge distillation (Kim and Rush, 2016, KD). We therefore investigate whether this can also reduce flicker in simultaneous translation. Since the initial model and the distilled model have the same architecture in our work, approximating KD is essentially self-training. We show that a self-trained model is able to achieve the same quality as the initial model, but with improved flicker-latency tradeoff. Furthermore, we show experiments that link the improved flicker to student monotonicity.

## 2 Background

### 2.1 Retranslation

We assume a retranslation approach, where the source is retranslated each time it is updated, and the new output replaces the old. Only the current sentence is retranslated – previous sentences are considered to be fixed. Retranslation can use an unmodified inference engine, in contrast to streaming approaches (e.g. (Ma et al., 2019a; Arivazhagan et al., 2019b)), making it simpler to deploy. The basic retranslation approach can be improved by using *prefix training* (Niehues et al., 2016, 2018), *biased beam search* and *output masking*<sup>2</sup> (Arivazhagan et al., 2020a).

<sup>1</sup>An example of this is shown in Appendix C

<sup>2</sup>This means that the last  $k$  words are omitted from the output before being passed to the user. It reduces flicker, but increases latency.

## 2.2 Evaluation of Simultaneous Translation

Evaluation of simultaneous translation requires that, as well as quality, we consider latency, and flicker (if we are using retranslation). The quality of the translation can be evaluated by comparing the final output of each sentence with a reference – we will use BLEU (Papineni et al., 2002; Post, 2018), CHRF (Popović, 2015) and COMET (Rei et al., 2020) scores. To measure flicker, we use *normalised erasure* (Arivazhagan et al., 2020a, 2019a), which measures the flicker between consecutive translation outputs by counting the minimum number of tokens that must be deleted from the end of the previous translation in order to produce the next, normalised by output length.

The measurement of latency has been the subject of some debate in the literature, with several different measures proposed (Ma et al., 2019a; Cherry and Foster, 2019; Ansari et al., 2021). In our experiments, we plot the flicker-latency tradeoff by controlling the output mask, and observing the effect on flicker. Since mask correlates with latency, our aim is to improve this mask-flicker tradeoff curve, and so be able to use a shorter mask with the same flicker budget.

## 2.3 Knowledge Distillation and Self-Training

The idea of sequence-level KD (Kim and Rush, 2016), is to create a smaller *student* model using the predictions of the larger *teacher* model. This has found application in MT efficiency (Junczys-Dowmunt et al., 2018) and in non-autoregressive translation (Zhou et al., 2020). For our purposes, the student model has the same size as the teacher. The output distributions of the self-trained model have lower entropy (Zhou et al., 2020), so the model is less likely to swap between translation hypotheses unnecessarily as the source prefix is extended. Also, since the self-trained model is trained on MT output, where the target order tends to be more similar to the source order (Zhou et al., 2020), it is more likely to avoid unnecessary reorderings, generating a more monotonic translation, which can be built up incrementally. We give experimental evidence for these in the next section.

Chen et al. (2021) also proposed to use pseudo-reference sentences obtained through forward translation of the source sentences to improve simultaneous translation. Unlike our work, they considered a streaming approach (specifically wait- $k$

(Ma et al., 2019b)) where the system can only append to the output; it does not flicker like retranslation. They showed that their approach could improve the quality-latency tradeoff of wait- $k$  using their distillation approach, but to create the training data for the student system they used wait- $k$  and filtering – we avoid these complications by just using the baseline system as the teacher.

## 3 Experiments

### 3.1 Data

We test our self-training approach on English $\leftrightarrow$ {German,Czech}. For En $\leftrightarrow$ De we use IWSLT21 (Anastasopoulos et al., 2021) for training, and the concatenation of the 2014 and 2015 test sets for development (early stopping), removing any sentences overlapping with training. For En $\leftrightarrow$ Cs, we use the training and validation set from WMT21 (Akhbardeh et al., 2021). Training data sizes are shown in Appendix A. We use *prefix training* to reduce the mismatch between sentence-level training data and prefix-based inference at test time (Niehues et al., 2018). For each parallel sentence pair in the training set, we generate a corresponding prefix pair by truncating using a randomly chosen proportion. We treat the validation sets similarly.

We test our systems both on IWSLT test data (derived from TED talks) and on the ESIC test set<sup>3</sup> (Macháek et al., 2021). From IWSLT, we use tst2018 for De $\leftrightarrow$ En, and tst2015/tst2016 combined for Cs $\leftrightarrow$ En. ESIC is derived from the European parliament proceedings, and consists of transcribed speeches in English, together with their simultaneous interpretation into Czech and German (also transcribed). ESIC is aligned at the document level, but not at the sentence level. We use the test portion for evaluation, only for En $\rightarrow$ X. It has been argued that such systems are better evaluated (and trained, if possible) on interpreted data (Zhao et al., 2021). However such data is hard to come by, and ESIC is the only such resource for European languages. We remove any segments from the IWSLT test sets that overlap with training, and also remove from the training data any europarl documents with overlap with ESIC.

All data is pre-processed with SentencePiece unigram model (Kudo and Richardson, 2018) with

<sup>3</sup><https://lindat.mff.cuni.cz/repository/xmlui/handle/11234/1-3719>

a shared subword (Sennrich et al., 2016b) vocabulary size of 32k.

Metric	Model	En→De		De→En		En→Cs		Cs→En	
		ESIC	IWSLT	IWSLT	ESIC	IWSLT	IWSLT	ESIC	IWSLT
BLEU	T	17.5	27.7	33.4	14.4	24.6	31.3		
	S	17.6	27.5	31.7	14.5	25.0	31.3		
ChrF	T	58.9	56.9	59.2	51.5	51.5	56.1		
	S	58.8	57.2	58.3	51.7	51.7	56.2		
COMET	T	.553	.330	.488	.651	.639	.519		
	S	.532	.326	.468	.672	.642	.521		

Table 1: Comparison between teacher (T) and student (S) models on ESIC and IWSLT test sets. For ESIC, BLEU and CHRf are calculated at document level, i.e. considering each document as a segment. For COMET we use reference-less `wmt20-comet-da` for ESIC and reference-based `wmt20-comet-da` for IWSLT.

### 3.2 Teacher-Student Training

Our teacher model, which serves as a baseline, is a transformer base (Vaswani et al., 2017) trained<sup>4</sup> with fairseq<sup>5</sup> (Ott et al., 2019).

We use the teacher to translate the training data, using a beam<sup>6</sup> of 8, then train a student model with the same architecture on this synthetic data.

In Table 1 we show the performance of our baseline system (equivalent to the teacher) and the student system on 6 test sets. Overall, student performance is robust compared to teacher, with same or better scores in Cs↔En and some small losses in De↔En.

To assess whether the student models reduce flicker in retranslation, we use each model in a simulated SLT pipeline and plot flicker-latency tradeoff curves. That is, we use the systems to translate ever-growing prefixes of the source sentences in the testsets, using SLTev (Ansari et al., 2021) to measure the flicker, and varying the output mask to show the tradeoff. A curve for one test set is shown in Figure 1, with full results in Appendix D. We can see that in all configurations the student models improve the flicker-latency tradeoff. In Appendix E, we show how the student training data is more monotonic, and the models have lower entropy, echoing Zhou et al. (2020).

<sup>4</sup>For training hyperparameters, see Appendix B.

<sup>5</sup>To generate training data for the students, we actually used a marian (Junczys-Dowmunt et al., 2018) model, with  $60 \times 10^6$  parameters, trained on the same data, with the same architecture, which achieves nearly identical BLEU. This was to take advantage of marian’s fast inference. All results shown in the paper are with the fairseq models.

<sup>6</sup>We also tested sequence-level interpolation, selecting the highest-scoring translation in an 8-best list according to BLEU and CHRf, but results were very similar.

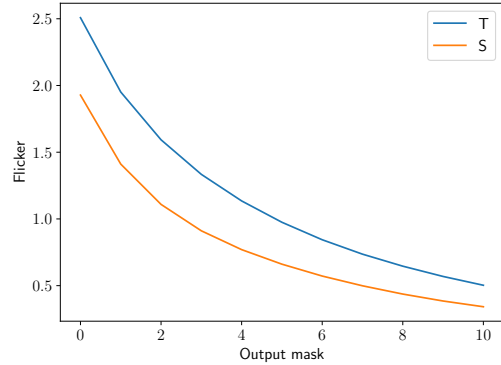


Figure 1: Flicker-latency tradeoff for the teacher (T) and student (S) models, En→De IWSLT. We control latency by varying the output mask.

### 3.3 Controlling Monotonicity

To show that self-training affects flicker through increased monotonicity, we experiment with controlling the monotonicity of the student training data. We stratify the teacher data into 5 different monotonicity levels using Kendall’s Tau on a *fast\_align* (Dyer et al., 2013) target–source alignment to measure monotonicity. We add the monotonicity level as pseudo-word, as in Sennrich et al. (2016a), to each source sentence, and train a teacher model on this monotonicity-aware corpus. We then use this teacher to create 5 different student training corpora, using the monotonicity control, and train 5 different students on these corpora.

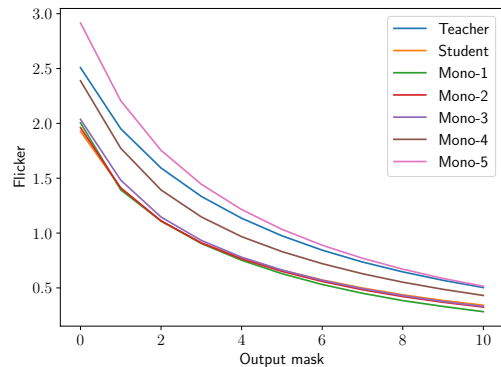


Figure 2: Latency–flicker tradeoff for the En→De IWSLT monotonicity-controlled models. Monotonicity control ranges from 1 (training data created with maximum monotonicity) to 5 (minimum monotonicity)

Table 2 shows the BLEU<sup>7</sup> scores for the monotonicity-controlled models, as well as the

<sup>7</sup>Scores for CHRf and COMET are in Appendix E, but the pattern is similar

teacher and student from the previous section. Using highly monotonic (Mono-1) or non-monotonic (Mono-5) data gives poor quality, but the in-between strata are similar, with Mono-3 slightly better overall. Figure 2 shows a distinctly worse flicker-latency tradeoff for Mono-5, whereas Mono-4 is a bit better than the teacher, and all other students are better. This supports the hypothesized connection between the higher degree of monotonicity in the student training data, and the better flicker-latency tradeoff in the student models.

Metric	Model	En→De		En→Cs	
		ESIC	IWSLT	ESIC	IWSLT
BLEU	Teacher	17.5	27.7	14.4	24.6
	Student	17.6	27.5	14.5	25.0
	Mono-1	8.6	14.4	14.7	23.6
	Mono-2	17.6	27.4	14.5	25.0
	Mono-3	17.5	27.9	14.5	25.7
	Mono-4	17.2	26.6	13.8	24.7
	Mono-5	16.0	25.0	12.5	23.0

Table 2: Student models with monotonicity control. Monotonicity ranges from 1 (highest) to 5 (lowest)

### 3.4 Self-training and Biased Beam Search

We investigate the combination of our self-training approach with biased beam search (Arivazhagan et al., 2020a). The idea of biased beam search (or “prefix biasing”) is to reduce flicker in retranslation by modifying inference so that the translation of the current prefix is “biased” towards the translation of the last prefix. At inference, the model has an extra term which penalises it for departing from the previous translation. As the current translation is being generated, once the hypothesis departs from the previous translation, we stop applying the bias, reverting to the unmodified MT model.

Before the previous translation is used for biasing, it is normally *masked*, i.e. the right-most  $k$  tokens are removed. Without applying this mask, biased beam search seriously reduces quality by forcing inference to follow poor quality early decisions. This *bias mask* is different from the output mask used in earlier experiments (which controls latency) although in previous work the bias and output mask are typically set to the same value.

We implemented biased beam search in fairseq and, based on previous work, we set the bias strength  $\beta = 0.25$ . After comparing different bias masks (Appendix F) we set the mask to 6 for ESIC and 10 for IWSLT.

We sweep across output masks to generate latency–flicker tradeoff curves in Figure 3 (with full results in Appendix F). We compare teacher and student models, with and without biased beam search. We can see from the graphs that biased beam search is effective in improving the latency–flicker tradeoff, but that the student models still improve over the teacher with biased beam search. The disadvantages of biased beam search are that it requires careful tuning of the prefix mask in order to avoid damaging quality, and that it requires a modified inference engine. The inference engine requires access to the previous translation, creating challenges for scalability. In contrast, our self-training approach requires no modifications to inference. Furthermore, since biased beam search relies on aligning the current translation with the previous one, it is hard to apply when the translation cannot be aligned – for example in a cascaded system where the ASR can rewrite its output.

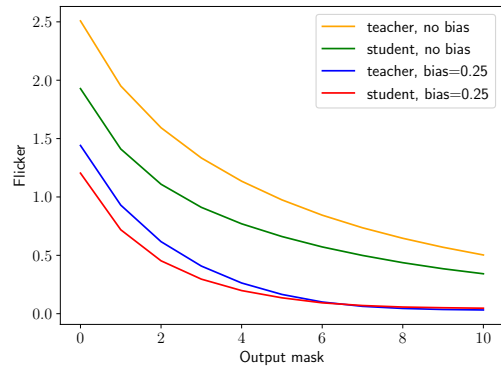


Figure 3: Latency-flicker tradeoff for teacher-student models with and without biased beam search, En→De, IWSLT

## 4 Conclusion

We show self-training reduces the flicker in retranslation-based simultaneous translation, whilst retaining quality. Our experiments link this flicker reduction to increased monotonicity and reduced entropy of the self-trained model. Although biased beam search can obtain larger reductions in flicker, it requires more careful parameter tuning, and a modified inference engine. However, one limitation of this work is that we evaluated it using only a couple of European language pairs (for which interpreted test sets are available). It will be interesting to see the results on language pairs with more syntactic divergence, in future.

293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
  
323  
324  
325  
326  
327  
328  
329  
  
330  
331  
332  
333  
334  
  
335  
336  
337  
338  
339  
340  
341  
  
342  
343  
344  
345  
346  
347  
348  
349  
  
350  
351

## References

Farhad Akhbardeh, Arkady Arkhangorodsky, Magdalena Biesialska, Ondej Bojar, Rajen Chatterjee, Vishrav Chaudhary, Marta R. Costa-jussa, Cristina España-Bonet, Angela Fan, Christian Federmann, Markus Freitag, Yvette Graham, Roman 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, Alahsera Auguste Tapo, Marco Turchi, Valentin Vydrin, and Marcos Zampieri. 2021. [Findings of the 2021 conference on machine translation \(WMT21\)](#). In *Proceedings of the Sixth Conference on Machine Translation*, pages 1–93, Online. Association for Computational Linguistics.

Antonios Anastasopoulos, Ondřej Bojar, Jacob Bremerman, Roldano Cattoni, Maha Elbayad, Marcello Federico, Xutai Ma, Satoshi Nakamura, Matteo Negri, Jan Niehues, Juan Pino, Elizabeth Salesky, Sebastian Stüker, Katsuhito Sudoh, Marco Turchi, Alexander Waibel, Changhan Wang, and Matthew Wiesner. 2021. [Findings of the IWSLT 2021 Evaluation Campaign](#). In *Proceedings of the 18th International Conference on Spoken Language Translation (IWSLT 2021)*, pages 1–29, Bangkok, Thailand (online). Association for Computational Linguistics.

Ebrahim Ansari, Ondřej Bojar, Barry Haddow, and Mohammad Mahmoudi. 2021. [SLTEV: Comprehensive evaluation of spoken language translation](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations*, pages 71–79, Online. Association for Computational Linguistics.

Naveen Arivazhagan, Colin Cherry, Te I, Wolfgang Macherey, Pallavi Baljekar, and George Foster. 2019a. [Re-Translation Strategies For Long Form, Simultaneous, Spoken Language Translation](#). *arXiv e-prints*, page arXiv:1912.03393.

Naveen Arivazhagan, Colin Cherry, Te I, Wolfgang Macherey, Pallavi Baljekar, and George Foster. 2020a. [Re-Translation Strategies For Long Form, Simultaneous, Spoken Language Translation](#). In *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*.

Naveen Arivazhagan, Colin Cherry, Wolfgang Macherey, Chung-Cheng Chiu, Semih Yavuz, Ruoming Pang, Wei Li, and Colin Raffel. 2019b. [Monotonic Infinite Lookback Attention for Simultaneous Machine Translation](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1313–1323, Florence, Italy. Association for Computational Linguistics.

Naveen Arivazhagan, Colin Cherry, Wolfgang Macherey, and George Foster. 2020b. [Re-translation](#)

[versus streaming for simultaneous translation](#). In *Proceedings of the 17th International Conference on Spoken Language Translation*, pages 220–227, Online. Association for Computational Linguistics.

Junkun Chen, Renjie Zheng, Atsuhito Kita, Mingbo Ma, and Liang Huang. 2021. [Improving simultaneous translation by incorporating pseudo-references with fewer reorderings](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5857–5864, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Colin Cherry and George Foster. 2019. [Thinking Slow about Latency Evaluation for Simultaneous Machine Translation](#). *arXiv e-prints*, page arXiv:1906.00048.

Chris Dyer, Victor Chahuneau, and Noah A. Smith. 2013. [A simple, fast, and effective reparameterization of IBM model 2](#). In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 644–648, Atlanta, Georgia. Association for Computational Linguistics.

Jiatao Gu, James Bradbury, Caiming Xiong, Victor O. K. Li, and Richard Socher. 2018. [Non-Autoregressive Neural Machine Translation](#). In *Proceedings of ICLR*.

Marcin Junczys-Dowmunt, Kenneth Heafield, Hieu Hoang, Roman Grundkiewicz, and Anthony Aue. 2018. [Marian: Cost-effective High-Quality Neural Machine Translation in C++](#). In *Proceedings of WNMT*.

Yoon Kim and Alexander M. Rush. 2016. [Sequence-Level Knowledge Distillation](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1317–1327. Association for Computational Linguistics. Event-place: Austin, Texas.

Taku Kudo and John Richardson. 2018. [SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71, Brussels, Belgium. Association for Computational Linguistics.

Mingbo Ma, Liang Huang, Hao Xiong, Renjie Zheng, Kaibo Liu, Baigong Zheng, Chuanqiang Zhang, Zhongjun He, Hairong Liu, Xing Li, Hua Wu, and Haifeng Wang. 2019a. [STACL: Simultaneous Translation with Implicit Anticipation and Controllable Latency using Prefix-to-Prefix Framework](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3025–3036, Florence, Italy. Association for Computational Linguistics.

407	Mingbo Ma, Liang Huang, Hao Xiong, Renjie Zheng,	<i>Human Language Technologies</i> , pages 35–40, San	464
408	Kaibo Liu, Baigong Zheng, Chuanqiang Zhang,	Diego, California. Association for Computational	465
409	Zhongjun He, Hairong Liu, Xing Li, Hua Wu, and	Linguistics.	466
410	Haifeng Wang. 2019b. <a href="#">STACL: Simultaneous trans-</a>		
411	<a href="#">lation with implicit anticipation and controllable la-</a>		
412	<a href="#">tency using prefix-to-prefix framework</a> . In <i>Proceed-</i>		
413	<i>ings of the 57th Annual Meeting of the Association</i>	Rico Sennrich, Barry Haddow, and Alexandra Birch.	467
414	<i>for Computational Linguistics</i> , pages 3025–3036,	2016b. <a href="#">Neural machine translation of rare words</a>	468
415	Florence, Italy. Association for Computational Lin-	<a href="#">with subword units</a> . In <i>Proceedings of the 54th An-</i>	469
416	guistics.	<i>annual Meeting of the Association for Computational</i>	470
		<i>Linguistics (Volume 1: Long Papers)</i> , pages 1715–	471
417	Dominik Macháek, Matú ilinec, and Ondej Bojar.	1725, Berlin, Germany. Association for Computa-	472
418	2021. <a href="#">Lost in Interpreting: Speech Translation from</a>	tional Linguistics.	473
419	<a href="#">Source or Interpreter?</a> In <i>Proceedings of Inter-</i>		
420	<i>speech</i> .	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	474
		Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz	475
421	Jan Niehues, Thai Son Nguyen, Eunah Cho, Thanh-Le	Kaiser, and Illia Polosukhin. 2017. <a href="#">Attention Is All</a>	476
422	Ha, Kevin Kilgour, Markus Müller, Matthias Sper-	<a href="#">You Need</a> . <i>CoRR</i> , abs/1706.03762.	477
423	ber, Sebastian Stüker, and Alex Waibel. 2016. <a href="#">Dy-</a>		
424	<a href="#">namic Transcription for Low-latency Speech Trans-</a>	Jinming Zhao, Philip Arthur, Gholamreza Haffari,	478
425	<a href="#">lation</a> . In <i>Proceedings of Interspeech</i> .	Trevor Cohn, and Ehsan Shareghi. 2021. <a href="#">It is not as</a>	479
		<a href="#">good as you think! evaluating simultaneous machine</a>	480
426	Jan Niehues, Ngoc-Quan Pham, Thanh-Le Ha,	<a href="#">translation on interpretation data</a> . In <i>Proceedings of</i>	481
427	Matthias Sperber, and Alex Waibel. 2018. <a href="#">Low-</a>	<i>the 2021 Conference on Empirical Methods in Natu-</i>	482
428	<a href="#">latency neural speech translation</a> . In <i>Proceedings of</i>	<i>ral Language Processing</i> , pages 6707–6715, Online	483
429	<i>Interspeech</i> .	and Punta Cana, Dominican Republic. Association	484
		for Computational Linguistics.	485
430	Myle Ott, Sergey Edunov, Alexei Baevski, Angela	Chunting Zhou, Jiatao Gu, and Graham Neubig.	486
431	Fan, Sam Gross, Nathan Ng, David Grangier, and	2020. <a href="#">Understanding knowledge distillation in non-</a>	487
432	Michael Auli. 2019. <a href="#">fairseq: A fast, extensible</a>	<a href="#">autoregressive machine translation</a> . In <i>8th Inter-</i>	488
433	<a href="#">toolkit for sequence modeling</a> . In <i>Proceedings of</i>	<i>national Conference on Learning Representations,</i>	489
434	<i>NAACL-HLT 2019: Demonstrations</i> .	<i>ICLR 2020, Addis Ababa, Ethiopia, April 26-30,</i>	490
		<i>2020</i> . OpenReview.net.	491
435	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	<b>A Training Data</b>	492
436	Jing Zhu. 2002. <a href="#">Bleu: a Method for Automatic Eval-</a>		
437	<a href="#">uation of Machine Translation</a> . In <i>Proceedings of</i>	Corpus	Sentence pairs
438	<i>40th Annual Meeting of the Association for Com-</i>	English-German	
439	<i>putational Linguistics</i> , pages 311–318, Philadelphia,	Europarl	1.79 M
440	Pennsylvania, USA. Association for Computational	Rapid	1.45 M
441	Linguistics. Type: Conference proceedings (arti-	News Commentary	0.35 M
442	cle).	OpenSubtitle	22.51 M
		TED corpus	206 K
443	Maja Popović. 2015. <a href="#">chrF: character n-gram F-score</a>	MuST-C.v2	248 K
444	<a href="#">for automatic MT evaluation</a> . In <i>Proceedings of the</i>	English-Czech	
445	<i>Tenth Workshop on Statistical Machine Translation,</i>	Europarl	645 K
446	pages 392–395, Lisbon, Portugal. Association for	ParaCrawl	14 M
447	Computational Linguistics.	CommonCrawl	161 K
		News Commentary	260 K
448	Matt Post. 2018. <a href="#">A call for clarity in reporting BLEU</a>	CzEng2.0	36 M <sup>8</sup>
449	<a href="#">scores</a> . In <i>Proceedings of the Third Conference on</i>	Wikitles	410 K
450	<i>Machine Translation: Research Papers</i> , pages 186–	Rapid	452 K
451	191, Brussels, Belgium. Association for Computa-		
452	tional Linguistics.	<b>B Training Parameters</b>	494
453	Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon	The non-default hyperparameters for Fairseq are	495
454	Lavie. 2020. <a href="#">COMET: A neural framework for MT</a>	shown in Table 3.	496
455	<a href="#">evaluation</a> . In <i>Proceedings of the 2020 Conference</i>		
456	<i>on Empirical Methods in Natural Language Process-</i>	<b>C Example of Flicker</b>	497
457	<i>ing (EMNLP)</i> , pages 2685–2702, Online. Associa-		
458	tion for Computational Linguistics.	An example of a translation which flickers be-	498
		tween two similar possibilities is shown in Table	499
459	Rico Sennrich, Barry Haddow, and Alexandra Birch.		
460	2016a. <a href="#">Controlling politeness in neural machine</a>		
461	<a href="#">translation via side constraints</a> . In <i>Proceedings of</i>		
462	<i>the 2016 Conference of the North American Chap-</i>		
463	<i>ter of the Association for Computational Linguistics:</i>		

Param	Value
label-smoothing	0.1
criterion	label_smoothed_cross_entropy
patience	10
arch	transformer
optimizer	adam
adam-betas	0.9, 0.98
lr	5e-4
lr-scheduler	inverse_sqrt
warmup-updates	4000
clip-norm	0.0
weight-decay	0.0001
dropout	0.3
update-freq	2
max-tokens	3000
best-checkpoint-metric	bleu
maximize-best-checkpoint-metric	True

Table 3: Fairseq training hyperparameters (non-default) for 4 GPU training.

4.

## D Flicker-Latency Tradeoff

In Figure 4, we show the flicker-latency tradeoff for all language-pair and testset combinations.

## E Monotonicity and Entropy of Student Models

We claimed that student models have lower flicker because they produce more monotonic translations, with less unnecessary variation. Here we provide evidence to support those claims.

**Training data for student models is more monotonic** In order to calculate the monotonicity of the training data, we use Kendall’s tau distance. We first extract word alignments from the training data using *fast\_align* (Dyer et al., 2013) to forward-align source and target. For each sentence pair we express the alignment as a function  $a : i \rightarrow j$ , and construct the two lists  $1, \dots, T$  and  $a(1), \dots, a(T)$  where  $T$  is the target length. We then calculate the Kendall’s tau between the two lists, repeat for each sentence pair in the corpus, and average. We repeat the calculation for the original training data and for the student training set. The results are shown in Table 5. We can see that in all cases, the student training data is more monotonic than the original teacher training data.

**Student models have lower entropy distributions** For each of our models, we calculate the

mean per-token entropy, by considering the probability distribution over the vocabulary at each time step. The entropies are shown in Table 6.

We can see from Table 6 that the token entropies are consistently lower for student models, showing that the distributions are more “peaky”, and so less likely to flicker between multiple output tokens with similar probabilities.

## F Biased beam search

We consider the effect of the bias mask on full sentence translation quality, as measured by BLEU. The bias mask is measured in sentencepiece tokens. Based on Figure 6 we set the bias mask to 6 for ESIC and 10 for IWSLT, in order to avoid a loss of BLEU.

In Figure 7, we show flicker-latency tradeoffs for all language pair and testset combinations.

<i>Source</i>	I hope you will have a little time and energy to focus on another report which is, despite its technicality, quite important for all of us.
<i>Target:</i>	Ich Ich hoffe, Ich hoffe, Sie Ich hoffe, Sie Ich hoffe, Sie haben Ich hoffe, Sie haben ein Ich hoffe, Sie werden ein wenig Zeit Ich hoffe, Sie haben etwas Zeit Ich hoffe, Sie haben etwas Zeit und Ich hoffe, Sie werden etwas Zeit und Energie haben, Ich hoffe, Sie haben etwas Zeit und Energie, um sich Ich hoffe, Sie haben etwas Zeit und Energie, um sich auf Ich hoffe, Sie werden ein wenig Zeit und Energie haben, um sich auf ein anderes Thema Ich hoffe, Sie haben etwas Zeit und Energie, um sich auf einen weiteren Bericht zu konzentrieren, Ich hoffe, Sie haben etwas Zeit und Energie, um sich auf einen anderen Bericht zu konzentrieren, ⋮ Ich hoffe, Sie werden ein wenig Zeit und Energie haben, um sich auf einen anderen Bericht zu konzentrieren, der trotz seiner Formalität für uns alle sehr wichtig ist.

Table 4: Examples of flicker caused by the teacher model. *Source* is the original full sentence which is input as a growing input prefix. *Target* is the output prefix in successive retranslations.

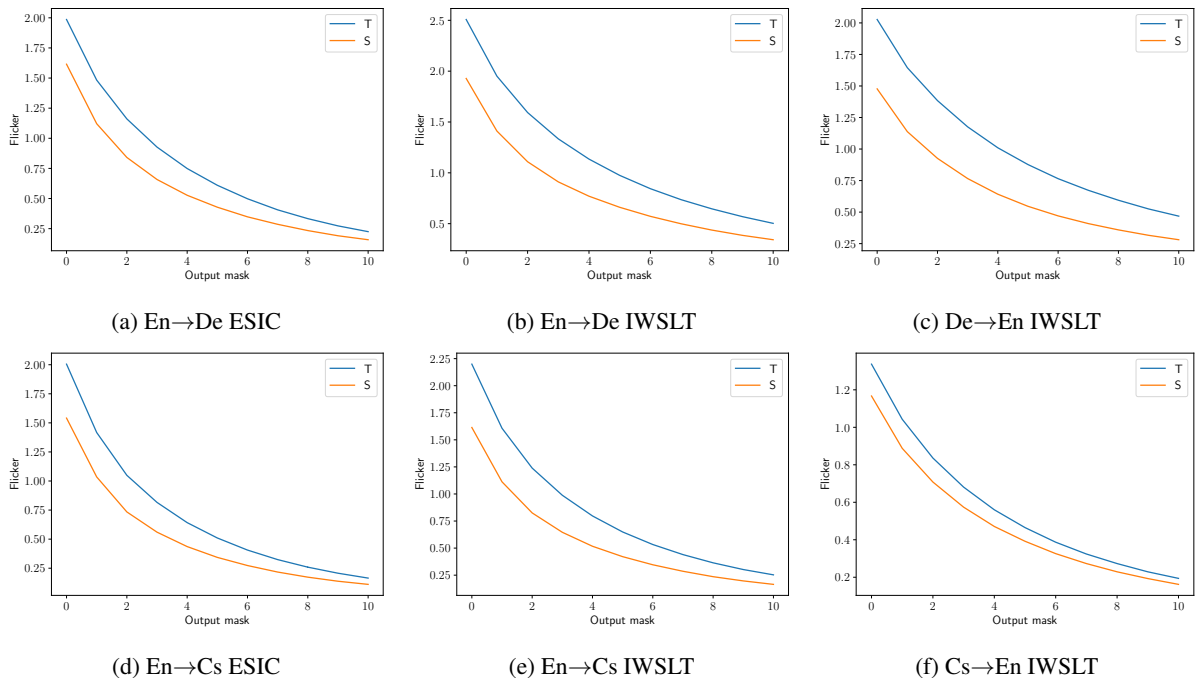


Figure 4: Flicker-latency tradeoff for the teacher-student models. We control latency by varying the output mask

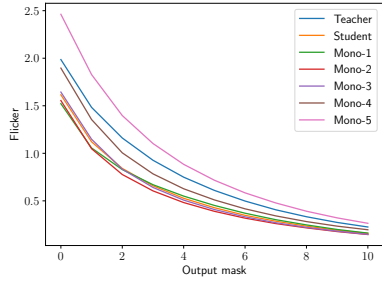
Model	En→De	De→En	En→Cs	Cs→En
Teacher	0.793	0.788	0.849	0.8356
Student	0.857	0.801	0.906	0.880

Table 5: Kendall’s tau distances. Higher scores indicate more monotonicity.

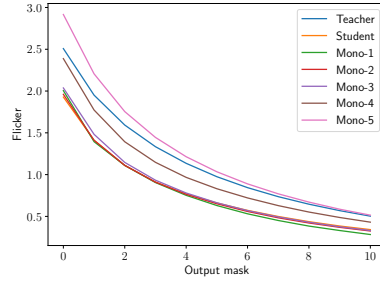
Pair	Test set	Entropy	
		Teacher	Student
En→De	ESIC	0.371	0.220
	IWSLT	0.295	0.228
De→En	IWSLT	0.273	0.160
	ESIC	0.443	0.251
En→Cs	IWSLT	0.417	0.238
	IWSLT	0.335	0.213

Table 6: Mean per-token entropies for each language pair test set combination

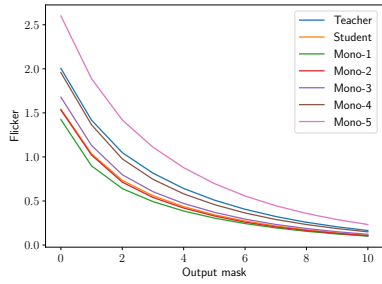




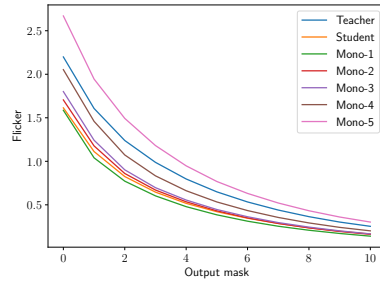
(a) En→De ESIC



(b) En→De IWSLT

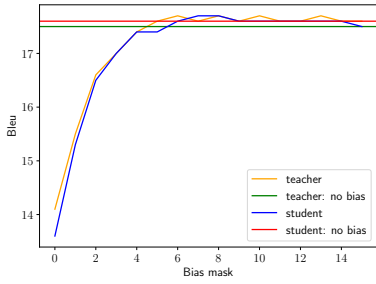


(c) En→Cs ESIC

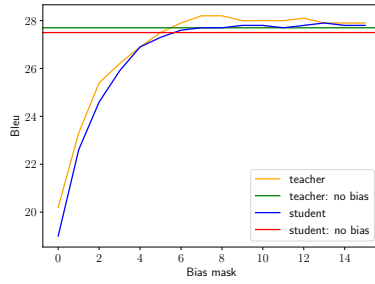


(d) En→Cs IWSLT

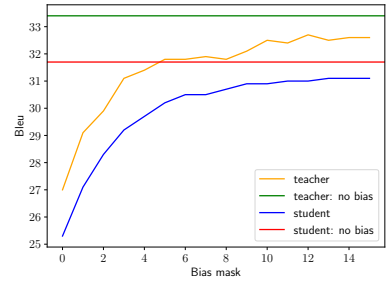
Figure 5: Latency–flicker tradeoff for the monotonicity-controlled models. Monotonicity control ranges from 1 (training data created with maximum monotonicity) to 5 (minimum monotonicity)



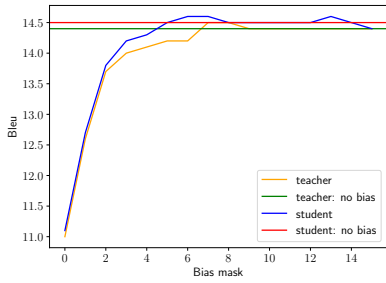
(a) En→De ESIC



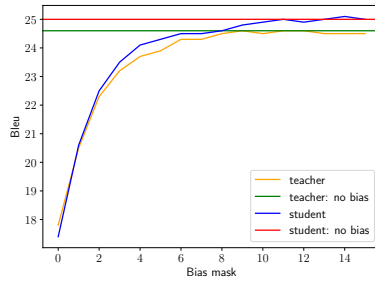
(b) En→De IWSLT



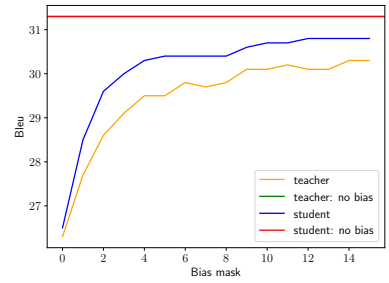
(c) De→En IWSLT



(d) En→Cs ESIC



(e) En→Cs IWSLT



(f) Cs→En IWSLT

Figure 6: Dependence of BLEU on bias mask when applying biased beam search.

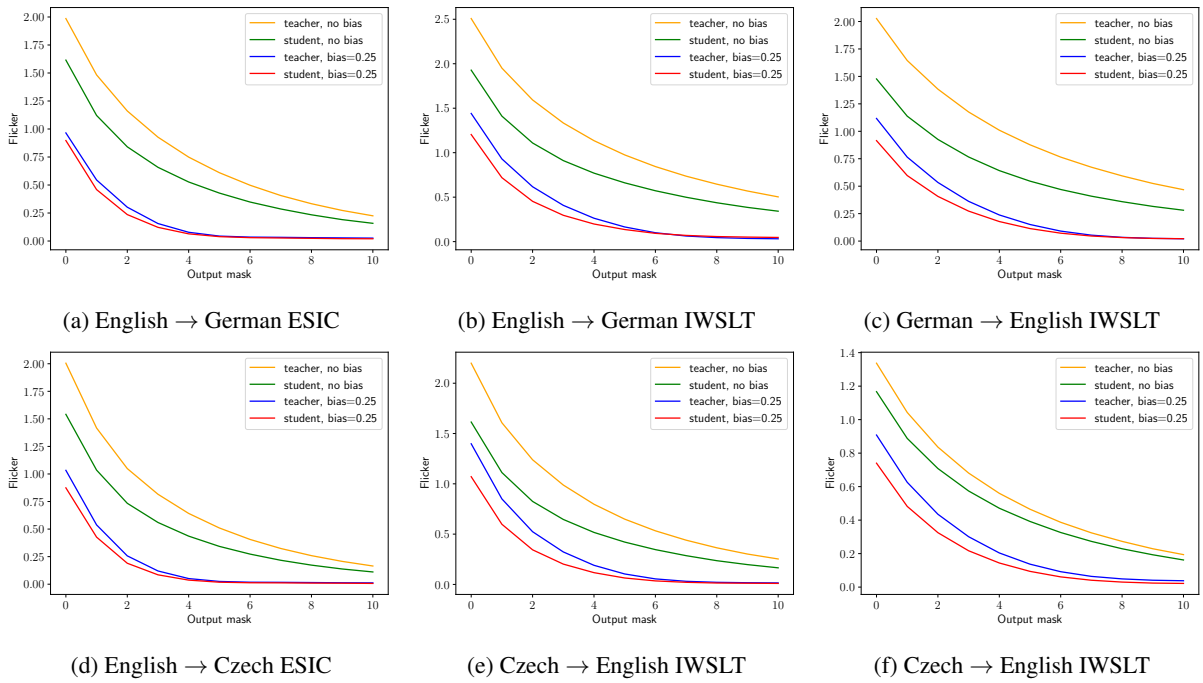


Figure 7: Flicker vs mask on biased beam search

Metric	Model	En→De		En→Cs	
		ESIC	IWSLT	ESIC	IWSLT
BLEU	Teacher	17.5	27.7	14.4	24.6
	Student <sub>model</sub>	17.6	27.5	14.5	25.0
	Mono-1	8.6	14.4	14.7	23.6
	Mono-2	17.6	27.4	14.5	25.0
	Mono-3	17.5	27.9	14.5	25.7
	Mono-4	17.2	26.6	13.8	24.7
	Mono-5	16.0	25.0	12.5	23.0
ChrF	Teacher	58.9	56.9	51.5	51.5
	Student <sub>model</sub>	58.8	57.2	51.7	51.7
	Mono-1	42.4	39.6	51.3	50.7
	Mono-2	58.7	57.3	51.8	52.0
	Mono-3	59.0	57.8	51.7	52.2
	Mono-4	59.0	56.8	51.4	51.4
	Mono-5	58.5	55.0	50.7	50.2
COMET	Teacher	.553	.330	.651	.639
	Student <sub>model</sub>	.532	.326	.672	.642
	Mono-1	.510	-0.028	.639	.597
	Mono-2	.526	.295	.650	.636
	Mono-3	.530	.326	.678	.641
	Mono-4	.535	.313	.677	.639
	Mono-5	.518	.247	.633	.577

Table 7: Full results of student models with monotonicity control.