

# STREAMING PUNCTUATION: A NOVEL PUNCTUATION TECHNIQUE LEVERAGING BIDIRECTIONAL CONTEXT FOR CONTINUOUS SPEECH RECOGNITION

Piyush Behre, Sharman Tan, Padma Varadharajan and Shuangyu Chang

Microsoft Corporation

## ABSTRACT

*While speech recognition Word Error Rate (WER) has reached human parity for English, continuous speech recognition scenarios such as voice typing and meeting transcriptions still suffer from segmentation and punctuation problems, resulting from irregular pausing patterns or slow speakers. Transformer sequence tagging models are effective at capturing long bi-directional context, which is crucial for automatic punctuation. Automatic Speech Recognition (ASR) production systems, however, are constrained by real-time requirements, making it hard to incorporate the right context when making punctuation decisions. Context within the segments produced by ASR decoders can be helpful but limiting in overall punctuation performance for a continuous speech session. In this paper, we propose a streaming approach for punctuation or re-punctuation of ASR output using dynamic decoding windows and measure its impact on punctuation and segmentation accuracy across scenarios. The new system tackles over-segmentation issues, improving segmentation  $F_{0.5}$ -score by 13.9%. Streaming punctuation achieves an average BLEU-score improvement of 0.66 for the downstream task of Machine Translation (MT).*

## KEYWORDS

*automatic punctuation, automatic speech recognition, re-punctuation, speech segmentation*

## 1. INTRODUCTION

Our hybrid Automatic Speech Recognition (ASR) generates punctuation with two systems working together. First, the decoder generates text segments and passes them to the Display Post Processor (DPP). The DPP system then applies punctuation to these text segments.

This two-stage setup works well for single-shot use cases such as voice assistant or voice search but performs poorly on long-form dictation. A dictation session typically comprises many spoken-form text segments generated by the decoder. Decoder features such as speaker pause duration determine the segment boundaries. The punctuation model in DPP then punctuates each of those segments. Without cross-segment look-ahead or the ability to correct previously finalized results, the punctuation model functions within the boundaries of each provided text segment. Consequently, punctuation model performance is highly dependent on the quality of text segments generated by the decoder.

Our past investments have focused on both systems independently - (1) improving decoder segmentation using look-ahead-based acoustic-linguistic features [32] and (2) using neural network architectures to punctuate in DPP. As measured by Punctuation- $F_1$  scores, these investments have improved our punctuation quality. However, over-segmentation in cases of slow speakers or irregular pausing is still prominent.

With streaming punctuation, we explore a system that discards decoder segmentation, instead shifting punctuation decision making towards a powerful long-context Transformer-based punctuation model. Rather than preliminary text segments, this system emits well-formed punctuated sentences, which is much more desirable for downstream tasks like translation of ASR output. The proposed architecture also satisfies real-time latency constraints for commercial ASR use cases.

Many works have demonstrated that leveraging prosodic features and audio inputs can improve punctuation quality [1, 2, 3]. However, as we show in our experiments, misleading pauses may significantly undermine punctuation quality and encourage overly aggressive punctuation. This is especially true in scenarios such as dictation, in which users pause often and unintentionally. Our work demonstrates that text-only streaming punctuation is robust to over-segmentation from irregular pauses and slow speakers.

We make the following key contributions: (1) We introduce a novel streaming punctuation approach to punctuate and re-punctuate ASR outputs, as described in section 3, (2) we demonstrate streaming punctuation's robustness to model architecture choices through experiments described in section 5, and (3) we achieve not only gains in punctuation quality but also significant downstream Bilingual Evaluation Understudy (BLEU) score gains on Machine Translation (MT) for a set of languages, as demonstrated in section 6.

## 2. RELATED WORK

Decoder segmentation conventionally involves using predefined silence timeouts or Voice Activity Detectors (VADs) to identify end of segment boundaries in text sequences. A separate punctuation system then applies punctuation and capitalization on these segments. Advancements in end-of-segment boundary detection include the addition of model-based techniques based on acoustic features [29, 30, 31] as well as acoustic-linguistic features [32]. Prior work has also explored end-to-end systems for end-of-segment boundary detection, jointly segmenting and decoding audio inputs with a focus on long-form ASR [33]. In this paper, we explore a system that does away with decoder segmentation boundaries and instead shifts punctuation decisions towards a powerful long-context Transformer-based punctuation model. As a baseline comparison, we use a conventional VAD-based system to generate decoder segments that we later feed into a downstream LSTM punctuation model.

Approaches to punctuation restoration have evolved to capture surrounding context more effectively. Early sequence labelling approaches for punctuation restoration used  $n$ -grams to capture context [4]. However, this simple approach becomes unscalable as  $n$  grows large, and does not generalize well to unseen data. This approach limits the amount of context that can be used in punctuation prediction.

Classical machine learning approaches such as conditional random fields (CRFs) [5, 6, 7], maximum entropy models [8], and hidden Markov models (HMMs) [9] model more complex features by leveraging manual feature engineering. This manual process is slow and cumbersome, and the quality of these features is dependent on feature engineers. These dependencies challenge the effectiveness of these classical approaches.

Neural approaches mostly displaced manual feature engineering, opting instead to learn more complex features through deep neural models. Recurrent neural networks (RNNs), specifically Gated Recurrent Units (GRUs) and bidirectional Long Short-Term Memory (LSTM) networks, have advanced natural language processing (NLP) and punctuation restoration by specifically modelling long-term dependencies in the text [10, 11, 12, 13, 14]. Prior works have also

successfully used LSTMs with CRF layers [15, 16]. Most recently, using Transformers [17] and especially pre-trained embeddings from models such as Bidirectional Encoder Representations from Transformers (BERT) [18] has significantly advanced quality across natural language processing (NLP) tasks. By leveraging attention mechanisms and more complex model architectures, Transformers can better capture bidirectional long-range text dependencies for punctuation restoration [19, 20, 21, 22, 23, 24].

### 3. PROPOSED METHOD

#### 3.1. Punctuation Model

We frame punctuation prediction as a neural sequence tagging problem. Figure 1 illustrates the end-to-end punctuation tagging and tag application process. We first tokenize the raw input text segment as a sequence of byte-pair encoding (BPE) tokens and pass this through a transformer encoder. Next, a punctuation token classification head, consisting of a dropout layer and a fully connected layer, generates token-level punctuation tags. Finally, we convert the token-level tags to word-level tags and generate the final punctuated text by appending each tag-specified punctuation symbol to the corresponding word in the input segment.

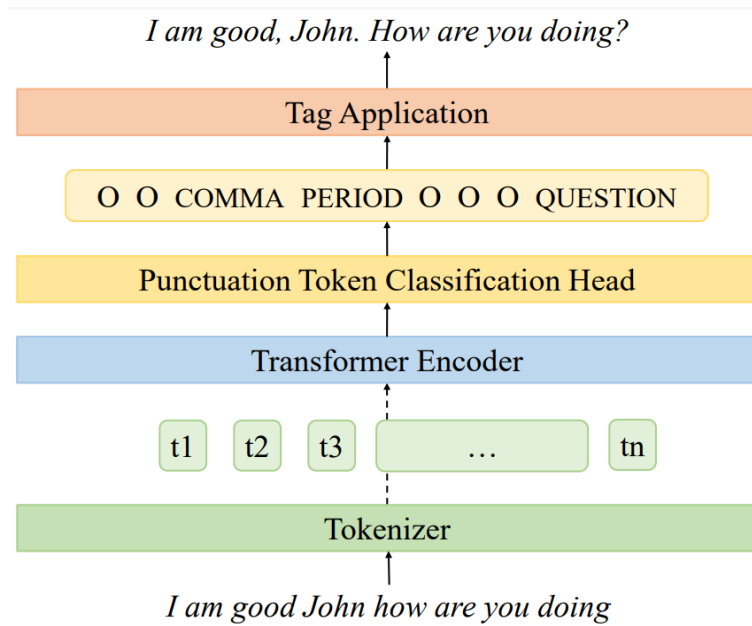


Figure 1. Punctuation tagging model using transformer encoder

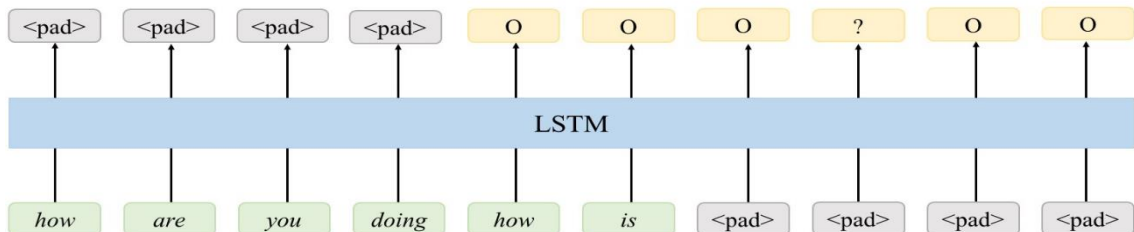


Figure 2. LSTM punctuation tagging model workflow with 4-word look-ahead enabled via <pad> input and output tokens

To demonstrate the robustness of streaming punctuation to different model architectures, we also conduct experiments with an LSTM tagging model with look-ahead rather than more powerful transformer-based models. Example text inputs to and tag outputs from the LSTM model with four-word look-ahead are illustrated in Figure 2 above. As the figure shows, we pad the model's inputs and outputs by tokens corresponding to specified look-ahead value to enable predictions with limited right context available. Once the model produces tags, we strip off the pad tokens to obtain a one-to-one mapping between each non-pad input token and each punctuation tag output.

### 3.2. Establishing Importance of Context Across Segment Boundaries

We first performed a study to better understand the importance and limitations of context as used in punctuating speech recognition output. Typically, ASR systems use silence-based timeouts or voice activity detection (VAD) to produce decoder segments. For slow speakers and users speaking with irregular pauses, this system can easily segment too aggressively. Similarly for fast speakers, the decoder segmentation may under-segment, resulting in lengthy segments. In an under-segmenting system, the segmentation eventually happens based on a pre-determined segment length timeout (e.g., 40 seconds).

Our baseline for this preliminary experiment was a system that punctuates only based on current context information. We considered two candidate systems for comparison. Left-segment context (LC) system considers previous segment appended as left context to the current segment but does not change the previous segment punctuation already produced. Right-segment context (RC) system considers the next segment appended as right context to the current segment, and only applies punctuation to the current segment. Table 1 presents the results of this preliminary experiment.

Table 1. Segmentation results with varying cross-segment context on a mixed set

Context setup	Segmentation					
	P	R	F <sub>1</sub>	F <sub>1</sub> -gain	F <sub>0.5</sub>	F <sub>0.5</sub> -gain
In-segment context	64	82	72		67	
Left-segment context	64	84	73	1.4%	67	0.4%
Right-segment context	80	69	74	2.8%	78	15.8%

Lower-precision and higher-recall is indicative of a system with over-segmentation problems. The LC system only slightly benefits from the additional left context and segmentation F<sub>0.5</sub> improves by only 0.4 percent. However, the RC system does a much better job in tackling the over-segmentation problem. This system improves segmentation F<sub>0.5</sub> by 15.8 percent, inverting the precision-recall tilt, which is much more desirable by our users. The RC system described here is not deployable in a streaming ASR service. However, this formed the basis of our streaming decoder approach to apply punctuation which we discuss next.

### 3.3. Streaming Decoder for Punctuation

Hybrid ASR systems often define segmentation boundaries using predefined silence thresholds. However, for human2machine scenarios like dictation, pauses do not necessarily indicate ideal segmentation boundaries for the ASR system. In our experience, users pause at unpredictable moments as they stop to think. All A1-4 segments in Table 2 are possible; each is a valid sentence with correct punctuation. Even with a punctuation model, if A4 is the user's intended sentence, all A1-3 would be incorrect. For dictation users, this system would produce over-

segmentation. To solve this issue, we must incorporate the right context across segment boundaries.

Table 2. Examples of possible segments generated by ASR

Id	Segment
Segment A1	It can happen.
Segment A2	It can happen in New York.
Segment A3	It can happen in New York City.
Segment A4	It can happen in New York City, right?

Our solution is a streaming punctuation system. The key is to emit complete sentences only after detecting the beginning of a new sentence. At each step, we punctuate text within a dynamic decoding window. This window consists of a buffer for which the system has not yet detected a sentence boundary as well as the new incoming segment. When the system detects at least one sentence boundary within the dynamic decoding window, we emit all complete sentences and reserve any remaining text as the new buffer. This process is illustrated in Figure 3 below.

Processing and punctuating each of the input segments separately and independently is problematic, as bad decoder segmentation leads to mid-sentence segmentation. As demonstrated in the punctuation finalized output column, streaming punctuation effectively enables punctuation decision-making across sentence boundaries with just enough surrounding context available. Streaming punctuation is a powerful way to improve final punctuation results, regardless of the decoder segmentation mechanism that is used to produce the input segments.

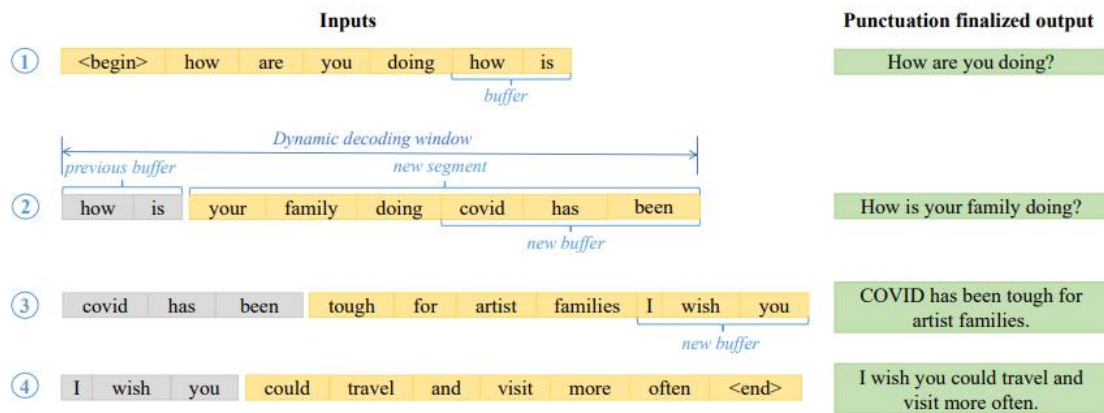


Figure 3. Dynamic decoding window for streaming punctuation

This strategy discards the original decoder boundary and decides the sentence boundary purely based on linguistic features. A powerful transformer model that captures the long context well is ideal for this strategy, as dynamic windows ensure that we incorporate enough left and right context before finalizing punctuation. Our approach also meets real-time requirements for ASR without incurring additional user-perceived latency, owing to the continual generation of hypothesis buffers within the same latency constraints. An improvement to this system would be to use a prosody-aware punctuation model that captures both acoustic and linguistic features. That would be a way to re-capture the acoustic cues that we lose by discarding the original segments. However, prosody-aware punctuation models may cause regressions in scenarios such as dictation in which users' pauses do not necessarily correspond to the presence of mid-sentence or end-of-sentence punctuation.

## 4. DATA PROCESSING PIPELINE

### 4.1. Datasets

We use public datasets from various domains to ensure a good mix of conversational and written-form data. Table 3 shows the word count distributions by percentage among the sets.

**OpenWebText** [25]: This dataset consists of web content from URLs shared on Reddit with at least three upvotes. This is our primary source of form written form (*human2machine*) data.

**Stack Exchange**: This dataset consists of user-contributed content on the Stack Exchange network. As this dataset consists of questions and answers, it is primarily of conversational (*human2human*) flavour.

**OpenSubtitles2016** [26]: This dataset consists of movie and TV subtitles. This is also primarily conversational (*human2human*).

**Multimodal Aligned Earnings Conference (MAEC)** [27]: This dataset consists of transcribed earnings calls based on S&P 1500 companies. Typically, each earnings call consists of a section of prepared remarks (*human2group*), followed by a Q&A section.

**National Public Radio (NPR) Podcast**: This dataset consists of transcribed NPR Podcast episodes. Typically, this consists of conversations between two to three individuals.

Table 3. Data distribution by number of words per training dataset

Dataset	Distribution
OpenWebText	52.8%
Stack Exchange	31.5%
OpenSubtitles2016	7.6%
MAEC	6.7%
NPR Podcast	1.4%

### 4.2. Data Processing

As described in Section 3.1, the transformer sequence tagging model takes spoken-form unpunctuated text as input and outputs a sequence of token-level tags signifying the punctuation to append to the corresponding input word. All datasets consist of punctuated written-form paragraphs, and we process them to generate spoken-form input text and corresponding output punctuation tag sequences for training.

To preserve the original context, we keep the original paragraph breaks in the datasets and use each paragraph as a training row. We first clean and filter the sets, removing symbols apart from alphanumeric, punctuation, and necessary mid-word symbols such as hyphens. To generate spoken-form unpunctuated data, we strip off all punctuation from the written-form paragraphs and use a Weighted Finite-State Transducers (WFST) based text normalization system to generate spoken-form paragraphs. During text normalization, we preserve alignments between each written-form word and its spoken form. We then use these alignments and the original punctuated display text to generate ground truth punctuation tags corresponding to the spoken-form text.

We set aside 10 percent or at most fifty thousand paragraphs from each set for validation and use the remaining data for training.

### 4.3. Tag Classes

We define four tag categories: comma, period, question mark, and ‘O’ for no punctuation. Each punctuation tag represents the punctuation symbol that appears appended to the corresponding text token. When we convert input word sequences into BPE sequences, we attach the tags only to the last BPE token for each word. We tag the rest of the tokens with ‘O’. For punctuation symbols other than comma, period, and question mark, we either convert them into the appropriate supported symbols (comma, period, or question mark) or remove those unsupported symbols entirely based on simple heuristics

## 5. EXPERIMENTS

### 5.1. Test Sets

We evaluate our punctuation model performance across various scenarios using private and public test sets. Each set contains long-form audio and corresponding written-form transcriptions with number formatting, capitalization, and punctuation. Starting from audio rather than text is critical to highlight the challenges associated with irregular pauses or slow speakers. This prohibits us from using the text-only International Conference on Spoken Language Translation (IWSLT) 2011 TED Talks corpus, typically used for reporting punctuation model performance.

**Dictation (Dict-100):** This internal set consists of one hundred sessions of long-form dictation ASR outputs and corresponding human transcriptions. On average, each session is 180 seconds long. Multiple judges process these sessions to generate the reference transcription in spoken and written form.

**MAEC:** 10 hours of test data taken from the MAEC corpus, containing transcribed earnings calls. This corresponds to ten earnings calls, each an hour long. Transcribers remove disfluencies, false starts, and repetitions for this set to make it more readable.

**European Parliament (EP-100):** This dataset contains one hundred English sessions scraped from European Parliament Plenary [34] videos. This dataset already contains English transcriptions, and human annotators provided corresponding translations into seven other languages. We use the source English transcriptions to measure segmentation and punctuation improvements. We use the translation reference to measure BLEU scores for the downstream task of Machine Translation.

**NPR Podcast (NPR-76):** 20 hours of test data from transcribed NPR Podcast episodes. On average, each session is 15 minutes long.

### 5.2. Experimental Setup

Our baseline system primarily uses Voice Activity Detection (VAD) based segmentation with a silence-based timeout threshold of 500ms. When VAD does not trigger, the system applies a segmentation at 40 seconds. The streaming punctuation system receives the input from the baseline system but can delay finalizing punctuation decisions until it detects the beginning of a new sentence.

We hypothesize that streaming punctuation outperforms the baseline system. We evaluate our hypothesis on LSTM and transformer punctuation tagging models. For the LSTM tagging model, we trained a 1-layer LSTM with 512-dimension word embeddings and 1024 hidden units. We used a look-ahead of four words, providing limited right context for better punctuation decisions. For the transformer tagging model, we trained a 12-layer transformer with sixteen attention heads, 1024-dimension word embeddings, 4096-dimension fully connected layers, and 8-dimension layers projecting from the transformer encoder to the decoder that maps to the tag classes.

We use 32 thousand BPE units for model input vocabulary. We limited training paragraph lengths to 250 BPE tokens and trimmed each to its last complete sentence. We trained all the models to convergence.

## 6. RESULTS AND DISCUSSION

We compare the results of our baseline (BL) and streaming (ST) punctuation systems on (1) the LSTM tagging model and (2) the Transformer tagging model. As expected, Transformers outperform LSTMs for this task. Here we evaluate our hypothesis for both model types to establish the effectiveness and robustness of our proposed system. For LSTM tagging models, BL-LSTM refers to the baseline system, and ST-LSTM refers to the streaming punctuation system. Similarly, for Transformer tagging models, BL-Transformer refers to the baseline system, and ST-Transformer refers to the streaming punctuation system.

### 6.1. Punctuation and Segmentation Accuracy

We measure and report punctuation accuracy with word-level precision (P), recall (R), and  $F_1$ -score. Table 4 summarizes punctuation metrics measured and aggregated over three punctuation categories: period, question mark, and comma.

Our customers consistently prefer higher precision (system only acting when confident) over higher recall (system punctuating generously). Punctuation- $F_1$  does not fully capture this preference. Customers also place higher importance on correctly detecting sentence boundaries over commas. We, therefore, propose segmentation- $F_{0.5}$  as a primary metric for this and future sentence segmentation work. The segmentation metric ignores commas and treats periods and question marks interchangeably, thus only measuring the quality of sentence boundaries. Table 5 summarizes segmentation metrics.

Although our target scenario was long-form dictation (human2machine), we found this technique equally beneficial for conversational (human2human) and broadcast (human2group) scenarios, establishing its robustness across applications. On average, the ST-Transformer system has a Segmentation- $F_{0.5}$  gain of 13.9 percent and a Punctuation- $F_1$  gain of 4.3 percent over the BL-Transformer system. Similarly, the ST-LSTM system has a Segmentation- $F_{0.5}$  improvement of 12.2 percent and a Punctuation- $F_1$  improvement of 2.1 percent over the BL-LSTM system. These results support our hypothesis that our streaming punctuation technique is effective and robust to different model architectures.

### 6.2. Downstream Task: Machine Translation

We measure the impact of segmentation and punctuation improvements on the downstream task of MT. Higher quality punctuation leads to translation BLEU gains for all seven target languages, as summarized in Table 6. The ST-Transformer system achieves the best results across all seven



target languages. On average, the ST-Transformer system has a BLEU score improvement of 0.66 over the BL-Transformer and wins for all target languages. Similarly, the ST-LSTM system has a BLEU score improvement of 0.33 over the BL-LSTM system and wins for five out of seven target languages. These results support our hypothesis.

Table 4. Punctuation results

Test Set	Model	PERIOD			Q-MARK			COMMA			OVERALL			F <sub>1</sub> -Gain
		P	R	F <sub>1</sub>	P	R	F <sub>1</sub>	P	R	F <sub>1</sub>	P	R	F <sub>1</sub>	
<i>Dict-100</i>	BL-LSTM	64	71	67	47	88	61	62	52	57	63	61	61	0.6%
	ST-LSTM	77	63	69	67	71	69	60	52	56	68	57	62	
	BL-Transf	69	76	72	50	88	64	68	52	59	68	63	65	2.9%
	ST-Transf	81	71	76	82	82	82	69	51	59	74	60	67	
<i>MAEC</i>	BL-LSTM	68	79	73	46	44	45	63	50	56	65	63	64	0.0%
	ST-LSTM	77	70	73	65	45	54	60	51	55	68	60	64	
	BL-Transf	71	80	75	50	50	50	65	49	56	67	63	65	2.4%
	ST-Transf	80	78	79	69	46	56	65	48	55	72	62	66	
<i>EP-100</i>	BL-LSTM	56	71	63	64	62	63	55	47	51	56	58	56	4.2%
	ST-LSTM	70	62	66	69	55	61	57	49	53	63	55	59	
	BL-Transf	58	76	66	58	70	64	57	49	53	57	61	59	5.8%
	ST-Transf	70	71	71	76	70	73	59	51	55	64	60	62	
<i>NPR-76</i>	BL-LSTM	72	71	72	71	66	69	65	58	61	69	65	67	4.0%
	ST-LSTM	82	71	76	76	69	73	65	59	62	74	66	70	
	BL-Transf	76	77	76	76	70	73	68	60	64	72	69	71	6.0%
	ST-Transf	87	79	83	81	75	78	70	61	65	79	71	75	

Table 5. Segmentation results

Test Set	Model	Segmentation					
		P	R	F <sub>1</sub>	F <sub>1</sub> -gain	F <sub>0.5</sub>	F <sub>0.5</sub> -gain
<i>Dict-100</i>	BL-LSTM	62	68	65		63	
	ST-LSTM	74	60	66	1.5%	71	12.0%
	BL-Transformer	66	74	70		67	
	ST-Transformer	79	69	73	4.3%	77	13.8%
<i>MAEC</i>	BL-LSTM	66	76	71		68	
	ST-LSTM	76	68	72	1.4%	74	9.5%
	BL-Transformer	69	77	73		70	
	ST-Transformer	79	75	77	5.5%	78	10.9%
<i>EP-100</i>	BL-LSTM	53	67	59		55	
	ST-LSTM	66	58	62	5.1%	64	16.1%
	BL-Transformer	54	72	62		57	
	ST-Transformer	67	68	68	9.7%	67	18.2%
<i>NPR-76</i>	BL-LSTM	71	70	70		71	
	ST-LSTM	81	70	75	7.1%	79	10.9%
	BL-Transformer	74	75	75		74	
	ST-Transformer	85	79	81	8.0%	84	12.5%

### 6.3. Downstream Task: Machine Translation

We measure the impact of segmentation and punctuation improvements on the downstream task of MT. Higher quality punctuation leads to translation BLEU gains for all seven target languages, as summarized in Table 6. The ST-Transformer system achieves the best results across all seven target languages. On average, the ST-Transformer system has a BLEU gain of 0.66 over BL-Transformer and wins for all target languages. Similarly, the ST-LSTM system has a BLEU gain of 0.33 over BL-LSTM system and wins for five out of seven target languages. These results support our hypothesis.

We used Azure Cognitive Services Translator API and compared them with reference translations. For Portuguese (pt) and French (fr), ST-LSTM regresses slightly, while ST-Transformer outperforms BL-Transformer. It is worth noting that ST-Transformer achieves significant gains over BL-Transformer, +1.1 for German (de) and +1.4 for Greek (el). The results suggest that punctuation has a higher impact on translation accuracy for some language pairs. For some language pairs, translation is more robust to punctuation errors.

Table 6. Translation BLEU Results: English audio recognized, punctuated, and translated to seven languages

Language	Model	BLEU	Gain
de	BL-LSTM	36.0	
	ST-LSTM	36.6	+0.6
	BL-Transformer	36.4	
	ST-Transformer	37.5	+1.1
el	BL-LSTM	39.8	
	ST-LSTM	40.8	+1.0
	BL-Transformer	40.3	
	ST-Transformer	41.7	+1.4
fr	BL-LSTM	41.0	
	ST-LSTM	40.6	-0.4
	BL-Transformer	41.7	
	ST-Transformer	41.8	+0.1
it	BL-LSTM	35.2	
	ST-LSTM	35.5	+0.3
	BL-Transformer	35.4	
	ST-Transformer	35.9	+0.5
pl	BL-LSTM	30.2	
	ST-LSTM	30.9	+0.7
	BL-Transformer	31.1	
	ST-Transformer	31.7	+0.6
pt	BL-LSTM	33.2	
	ST-LSTM	33	-0.2
	BL-Transformer	33.7	
	ST-Transformer	33.9	+0.2
ro	BL-LSTM	39.8	
	ST-LSTM	40.1	+0.3
	BL-Transformer	40.5	
	ST-Transformer	41.2	+0.7

Table 7. BL-Transformer's incorrect punctuation leads to incorrect translations from English to select four languages. ST-Transformer correctly punctuates, resulting in correct translations.

Language	BL-Transformer	ST-Transformer
en	I. Just have to share the view . . .	I just have to share the view . . .
de	I. Ich muss nur die Ansicht teilen . . .	Ich muss nur die Ansicht teilen . . .
fr	I. Il suffit de partager le point de . . .	Je dois simplement partager le point de . . .
it	I. Basti condividere l'opinione . . .	Devo solo condividere l'opinione . . .

Table 7 presents an example of how incorrect punctuation can lead to downstream consequences in machine translated outputs. Here BL-Transformer incorrectly punctuates after "I" which results in (1) failure to accurately translate the word, (2) incorrect translations for the subsequent text, and (3) incorrect punctuation in the translations to all languages. ST-Transformer, however, correctly punctuates and thus produces correct translations. This example demonstrates the importance of punctuation quality for downstream tasks such as MT.

## 7. CONCLUSION

Long pauses and hesitations occur naturally in dictation scenarios. We started this work to solve the over-segmentation problem for long-form dictation users. We discovered these elements affect other long-form transcription scenarios like conversations, meeting transcriptions, and broadcasts. Our streaming punctuation approach improves punctuation for a variety of these ASR scenarios. Higher quality punctuation directly leads to higher quality downstream tasks, such as improvement in BLEU scores for machine translation applications. We also established the efficacy of streaming punctuation across the transformer and LSTM tagging models, thus establishing the robustness of streaming punctuation to different model architectures.

In this paper, we focused on improving punctuation for hybrid ASR systems. Our preliminary analysis has found that though end-to-end (E2E) ASR systems produce better punctuation out of the box, such systems have yet to fully solve the problem of over-segmentation and could benefit from streaming re-punctuation techniques. We plan to present our findings in the future. Streaming punctuation discussed here relies primarily on linguistic features and discards acoustic signals. We plan to further extend this work using prosody-aware neural punctuation models. As we explore streaming punctuation's effectiveness and potential for other languages, we are also interested in exploring the impact of intonation or accents on our method.

## REFERENCES

- [1] Elizabeth Shriberg, Andreas Stolcke, Dilek HakkaniTur, and Gokhan Tur, "Prosody-based automatic segmentation of speech into sentences and topics," *Speech communication*, vol. 32, no. 1-2, pp. 127–154, 2000.
- [2] Madina Hasan, Rama Doddipatla, and Thomas Hain, "Multi-pass sentence-end detection of lecture speech," in *Fifteenth Annual Conference of the International Speech Communication Association*, 2014.
- [3] Piotr Zelasko, Piotr Szymański, Jan Mizgajski, Adrian Szymczak, Yishay Carmiel, and Najim Dehak, "Punctuation prediction model for conversational speech," arXiv preprint arXiv:1807.00543, 2018.
- [4] Agustin Gravano, Martin Jansche, and Michiel Bacchiani, "Restoring punctuation and capitalization in transcribed speech," in *2009 IEEE International Conference on Acoustics, Speech, and Signal Processing*. IEEE, 2009, pp. 4741–4744.
- [5] Wei Lu and Hwee Tou Ng, "Better punctuation prediction with dynamic conditional random fields," in *Proceedings of the 2010 conference on empirical methods in natural language processing*, 2010, pp. 177–186.
- [6] Xuancong Wang, Hwee Tou Ng, and Khe Chai Sim, "Dynamic conditional random fields for joint sentence boundary and punctuation prediction," in *Thirteenth Annual Conference of the International Speech Communication Association*, 2012.
- [7] Nicola Ueffing, Maximilian Bisani, and Paul Vozila, "Improved models for automatic punctuation prediction for spoken and written text.," in *Interspeech*, 2013, pp. 3097–3101.
- [8] Jing Huang and Geoffrey Zweig, "Maximum entropy model for punctuation annotation from speech.," in *Interspeech*, 2002.
- [9] Yang Liu, Elizabeth Shriberg, Andreas Stolcke, Dustin Hillard, Mari Ostendorf, and Mary Harper, "Enriching speech recognition with automatic detection of sentence boundaries and disfluencies," *IEEE Transactions on audio, speech, and language processing*, vol. 14, no. 5, pp. 1526–1540, 2006.
- [10] Ronan Collobert, Jason Weston, Leon Bottou, Michael, Karlen, Koray Kavukcuoglu, and Pavel Kuksa, "Natural language processing (almost) from scratch," *Journal of machine learning research*, vol. 12, no. ARTICLE, pp. 2493–2537, 2011.
- [11] Xiaoyin Che, Cheng Wang, Haojin Yang, and Christoph Meinel, "Punctuation prediction for unsegmented transcript based on word vector," in *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, 2016, pp. 654–658.

- [12] William Gale and Sarangarajan Parthasarathy, “Experiments in character-level neural network models for punctuation.,” in INTERSPEECH, 2017, pp. 2794–2798.
- [13] Vasile Pais, and Dan Tufis, “Capitalization and punctuation restoration: a survey,” *Artificial Intelligence Review*, vol. 55, no. 3, pp. 1681–1722, 2022.
- [14] Kaituo Xu, Lei Xie, and Kaisheng Yao, “Investigating lstm for punctuation prediction,” in 2016 10th International Symposium on Chinese Spoken Language Processing (ISCSLP), 2016, pp. 1–5.
- [15] Xuezhe Ma and Eduard Hovy, “End-to-end sequence labeling via bi-directional lstm-cnns-crf,” *arXiv preprint arXiv:1603.01354*, 2016.
- [16] Jiangyan Yi, Jianhua Tao, Zhengqi Wen, Ya Li, et al., “Distilling knowledge from an ensemble of models for punctuation prediction.,” in *Interspeech*, 2017, pp. 2779–2783.
- [17] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin, “Attention is all you need,” *Advances in neural information processing systems*, vol. 30, 2017.
- [18] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” *arXiv preprint arXiv:1810.04805*, 2018.
- [19] Jiangyan Yi, Jianhua Tao, Zhengkun Tian, Ye Bai, and Cunhang Fan, “Focal loss for punctuation prediction.,” in INTERSPEECH, 2020, pp. 721–725.
- [20] Yangjun Wu, Kebin Fang, and Yao Zhao, “A contextaware feature fusion framework for punctuation restoration,” *arXiv preprint arXiv:2203.12487*, 2022.
- [21] Maury Courtland, Adam Faulkner, and Gayle McElvain, “Efficient automatic punctuation restoration using bidirectional transformers with robust inference,” in *Proceedings of the 17th International Conference on Spoken Language Translation*, 2020, pp. 272–279.
- [22] Tanvirul Alam, Akib Khan, and Firoj Alam, “Punctuation restoration using transformer models for high-and low-resource languages,” in *Proceedings of the Sixth Workshop on Noisy User-generated Text (W-NUT 2020)*, 2020, pp. 132–142.
- [23] Raghavendra Pappagari, Piotr Zelasko, Agnieszka Mikołajczyk, Piotr Pezik, and Najim Dehak, “Joint prediction of truecasing and punctuation for conversational speech in low-resource scenarios,” in *2021 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2021, pp. 1185–1191.
- [24] Attila Nagy, Bence Bial, and Judit Acs, “Automatic punctuation restoration with bert models,” *arXiv preprint arXiv:2101.07343*, 2021.
- [25] Aaron Gokaslan and Vanya Cohen, “Openwebtext corpus,” <http://Skylion007.github.io/OpenWebTextCorpus>, 2019.
- [26] Pierre Lison and Jorg Tiedemann, “Opensubtitles2016: Extracting large parallel corpora from movie and tv subtitles,” 2016.
- [27] Jiazheng Li, Linyi Yang, Barry Smyth, and Ruihai Dong, “Maec: A multimodal aligned earnings conference call dataset for financial risk prediction,” in *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, 2020, pp. 3063–3070.
- [28] Piyush Behre, Sharman Tan, Padma Varadharajan and Shuangyu Chang, “Streaming Punctuation for Long-form Dictation with Transformers.” in *Proceedings of the 8<sup>th</sup> International Conference on Signal, Image Processing and Embedded Systems*, 2022, pp. 187-197. *arXiv:2210.05756*
- [29] Elizabeth Shriberg, Andreas Stolcke, Dilek Hakkani-Tur, and Gokhan Tur, “Prosody-based automatic segmentation of speech into sentences and topics,” *Speech communication*, vol. 32, no. 1-2, pp. 127–154, 2000.
- [30] Zulfiqar Ali and Muhammad Talha, “Innovative method for unsupervised voice activity detection and classification of audio segments,” *IEEE Access*, vol. 6, pp. 15494– 15504, 2018.
- [31] Junfeng Hou, Wu Guo, Yan Song, and Li-Rong Dai, “Segment boundary detection directed attention for online end-to-end speech recognition,” *EURASIP Journal on Audio, Speech, and Music Processing*, vol. 2020, no. 1, pp. 1–16, 2020.
- [32] Piyush Behre, Naveen Parihar, Sharman Tan, Amy Shah, Eva Sharma, Geoffrey Liu, Shuangyu Chang, Hosam Khalil, Chris Basoglu, and Sayan Pathak. “Smart Speech Segmentation using Acousto-Linguistic Features with look-ahead,” *arXiv preprint arXiv:2210.14446*, 2022.
- [33] W Ronny Huang, Shuo-yiin Chang, David Rybach, Rohit Prabhavalkar, Tara N Sainath, Cyril Allauzen, Cal Peyser, and Zhiyun Lu, “E2E Segmenter: Joint segmenting and decoding for long-form ASR,” *arXiv preprint arXiv:2204.10749*, 2022.
- [34] “Debates and videos: Plenary: European parliament” <https://www.europarl.europa.eu/plenary/en/debates-video.html>, Accessed: 2022-05-30.