

InfiniSST: Simultaneous Translation of Unbounded Speech with Large Language Model

Anonymous ACL submission

Abstract

Simultaneous translation of unbounded streaming speech remains a challenging problem due to the need for effectively processing the history speech context and past translations so that quality and latency, including computation overhead, can be balanced. Most prior works assume pre-segmented speech, limiting their real-world applicability. In this paper, we propose InfiniSST, a novel approach that formulates SST as a multi-turn dialogue task, enabling seamless translation of unbounded speech. We construct translation trajectories and robust segments from MuST-C with multi-latency augmentation during training and develop a key-value (KV) cache management strategy to facilitate efficient inference. Experiments on MuST-C En-Es, En-De, and En-Zh demonstrate that InfiniSST reduces computation-aware latency by 0.5 to 1 second while maintaining the same translation quality compared to baselines. Ablation studies further validate the contributions of our data construction and cache management strategy.

1 Introduction

Simultaneous speech translation (SST) is the task of translating partial speech input from a source language into text in a target language, with a wide range of applications, including conference interpretation and live-streaming translation (Ma et al., 2020b; Ren et al., 2020). Most prior research on SST focuses on translating pre-segmented speech (SST-S), assuming that gold-standard segmentation is provided (Liu et al., 2021; Zeng et al., 2021; Dong et al., 2022; Papi et al., 2023, 2024b). However, translating unbounded, streaming speech (SST-U) remains underexplored.

Unbounded speech presents a major challenge that the model has to effectively process the history speech context and past translations so that quality and latency, including computation overhead, can be balanced. Large language model (LLM)

is a promising solution for long-context modeling with the recent advancements (Su et al., 2021; Han et al., 2024). Moreover, LLM-based architectures have been shown to improve SST-S performance (Xu et al., 2024). However, conventional SST-S approaches suffer from high computational costs, as they require recomputing features for past speech and generated text every time a new speech chunk arrives. Some studies mitigate this issue by framing SST as a multi-turn dialogue task, either explicitly (Yu et al., 2025; Wang et al., 2024) or implicitly (Ouyang et al., 2024; Raffel et al., 2024), leveraging key-value (KV) caching to improve efficiency. While effective for segmented speech and text, these methods do not seamlessly extend to unbounded speech.

In this paper, we propose InfiniSST, a method for simultaneous translation of unbounded speech using a multi-turn dialogue format. We construct SST trajectories and derive robust speech segments for training from the MuST-C dataset, enhancing them with a multi-latency strategy to increase diversity. During inference, we employ a KV cache management strategy, inspired by Han et al. (2024), to enable seamless extrapolation to unbounded speech input. Experiments on MuST-C En-Es, En-De, and En-Zh (Di Gangi et al., 2019) show that InfiniSST reduces computation-aware latency by 0.5 to 1 second while maintaining the same BLEU score as baselines. A detailed ablation study further validates the effectiveness of our data construction and cache management strategies during inference.

2 Related Works

2.1 SST on Unbounded Speech

Cascade Approaches Cascade-based methods typically use an automatic speech recognition (ASR) model to segment and transcribe the input, followed by a machine translation model that

042
043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080

translates the transcription (Fugen et al., 2006; Yoshimura et al., 2020; Huang et al., 2022; Donato et al., 2021). However, segmentation errors and the lack of punctuation degrade translation quality, which complicates maintaining low latency and high quality.

Direct SST on Unbounded Speech Several works explore end-to-end approaches for SST on unbounded speech (Schneider and Waibel, 2020; Iranzo-Sánchez et al., 2024). These methods avoid external segmentation by dynamically preserving relevant audio context and previously generated text while discarding older information. Papi et al. (2024a) extends AlignAtt to unbounded speech by storing text and audio history in a fully streaming way, which helps reduce latency and maintain contextual awareness. Despite these advances, balancing translation quality, latency, and computational demands remains a challenge. Our approach addresses these issues by managing unbounded speech input without loss in translation accuracy and with improved computational efficiency.

2.2 Length Extrapolation of LLM

Recent advances in positional encoding (Su et al., 2021; Press et al., 2021; Sun et al., 2023) have enabled models to handle longer sequences with little or no additional training. ReRoPE (Su, 2023) introduces an NTK-aware Scaled RoPE that extends context length to infinite without fine-tuning. Han et al. (2024) and Xiao et al. (2024) propose on-the-fly length generalization based on a Λ -shaped attention window, allowing nearly unlimited input length with no fine-tuning. InfiniSST is a successful application of RoPE and Λ -shaped attention window in SST-U.

3 Method

3.1 Problem Formulation

Let $\mathbf{s}_{1:t} = (s_1, s_2, \dots, s_t)$ be the partial input of an unbounded input speech sequence and $\mathbf{y}_{1:i} = (y_1, y_2, \dots, y_i)$ represent the partial text translation. Here $\mathbf{s}_{1:t}$ is raw speech input instead of speech features. Define $\pi(\mathbf{s}_{1:t}, \mathbf{y}_{1:i}) \in [0, 1]$ as the policy to determine whether to take more speech input (=0) or to generate target translation tokens (=1). Whenever $\pi(\mathbf{s}_{1:t}, \mathbf{y}_{1:i}) = 1$, we define $g_{i+1} = t$ as the delay of $i + 1$ -th token. Let $g_0 = 0$. In addition, let θ be the model parameter, we define the probability of generating next token given a partial

speech input as $P_\theta(y_{i+1} | \mathbf{s}_{1:t}, \mathbf{y}_{1:i})$. In our formulation, we use a simple policy by checking whether the current generated token y_i is a special ending token T_0 (e.g. stop writing translation and read speech input when encountering “⟨EOT⟩” token in Llama (Grattafiori et al., 2024)).

$$\pi(\mathbf{s}_{1:t}, \mathbf{y}_{1:i}) = \begin{cases} 0, & \text{if } y_i = T_0 \\ 1, & \text{otherwise} \end{cases} \quad (1)$$

Given \mathbf{s} , we define the conditional probability of generating a translation sequence $\mathbf{y}_{1:i}$ with associated delays for each token $\mathbf{g}_{1:i}$ as:

$$P(\mathbf{y}, \mathbf{g} | \mathbf{s}) = \prod_{i=1}^{|\mathbf{y}|} \left(P_\theta(y_i | \mathbf{s}_{1:g_i}, \mathbf{y}_{1:i-1}) \right. \\ \left. \prod_{j=g_{i-1}}^{g_i-1} (1 - \pi(\mathbf{s}_{1:j}, \mathbf{y}_{i-1})) \right) \quad (2)$$

The translation quality and latency are subsequently evaluated based on \mathbf{s} , \mathbf{y} and \mathbf{g} .

3.2 Model Architecture

We design InfiniSST, a simultaneous speech translation model that can take unbounded streaming speech input and generate target text efficiently. The InfiniSST consists of 1) a streaming speech encoder to incrementally compute representations of partial speech input without recomputation, 2) a speech-to-token embedding adapter to match speech representations to LLM’s token embedding space, and 3) an multi-turn LLM decoder to interactively take speech input and generate translation as needed (Figure 1).

Streaming Speech Encoder We modify a pre-trained wav2vec2 (Baevski et al., 2020) speech encoder to encode the unbounded streaming speech input. However, there is major limitation of the original wav2vec2. It uses bidirectional attention and bidirectional convolutional position embedding, which needs to recompute the representations for every new segment of streaming speech input. To handle unbounded speech input, we introduce three modifications to the speech encoder. Firstly, we replace the wav2vec2’s convolutional positional embedding with a rotary positional embedding (RoPE) (Su et al., 2024) because it shows better extensibility for long sequences. Secondly, we replace bidirectional attention with chunk-wise causal attention (Deng et al., 2022). Each chunk

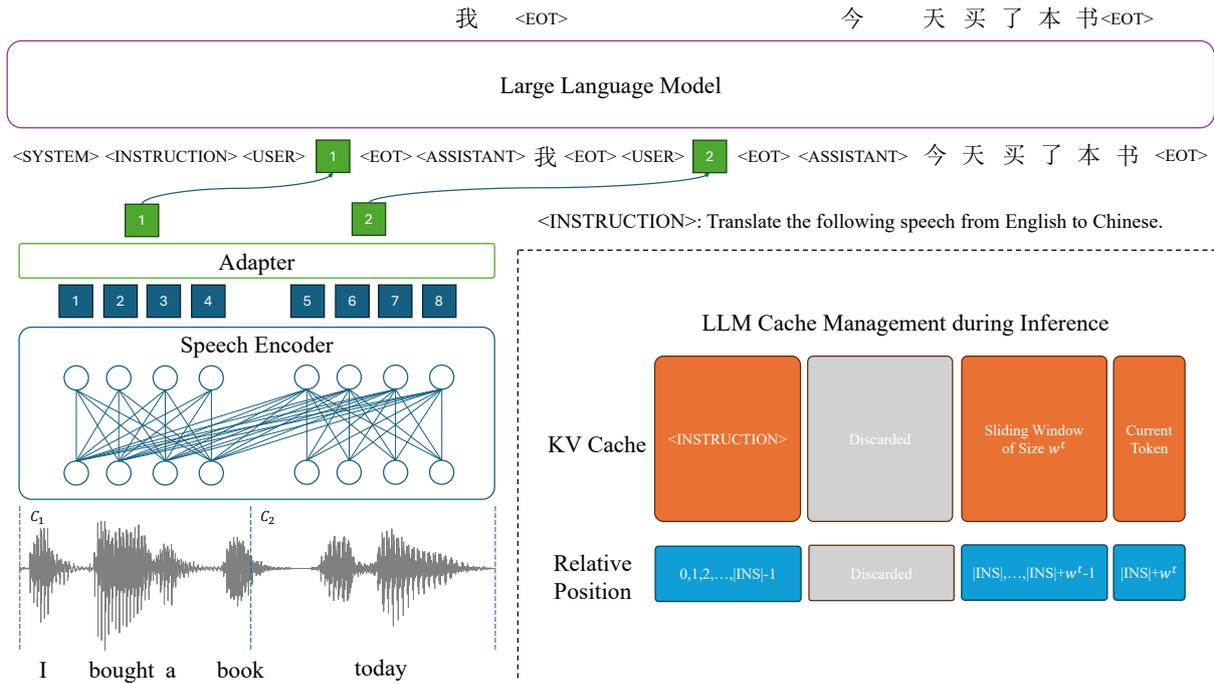


Figure 1: Model architecture of InfiniSST. InfiniSST first encodes speech using a chunkwise-causal speech encoder, then compresses the speech features into embeddings via an adapter. The large language model (LLM) processes the input by first reading a system instruction, then alternating between consuming speech embeddings and generating translations. The translation process stops when the LLM generates an EOT token. During inference, we employ a sliding window of size w^t for the LLM, conditioning the translation on the most recent w^t KV caches along with the KV cache of the system instruction, enabling extrapolation to unbounded speech input.

171 contains 48 frames in wav2vec2, with a total duration of 960ms. The multihead attention within each
 172 chunk remains bidirectional while attention across
 173 chunks is causal. This is achieved by adding block-
 174 wise masking to the attention weights. Thirdly, we
 175 apply a sliding window mechanism with window
 176 size w^s to maintain a finite context length, restrict-
 177 ing chunk i to attend only to hidden states of chunks
 178 $[i - w^s + 1, i]$. In practice, we use $w^s = 10$ so each
 179 speech embedding is computed from roughly 9.6
 180 seconds of the preceding speech input.
 181

182 **Speech-to-Token Embedding Adapter** The
 183 above speech encoder would produce a slightly
 184 longer sequence of embeddings compared to the
 185 lengths of corresponding transcript. The encoder’s
 186 output embeddings are also different from the to-
 187 ken embeddings of the later LLM. To reduce the
 188 length of the speech encoder output, we apply two
 189 1d convolutional layers with a kernel size of 2 and
 190 a stride of 2. We add a linear projection layer that
 191 maps from convolutional output to the LLM’s em-
 192 bedding space. Our speech-to-token embedding
 193 adapter downsamples the input by a factor of 4.
 194 Therefore a chunk with 48 frames of speech input

will result in 12 embedding vectors. 195

196 **Multi-turn LLM Decoder** Our decoder needs
 197 to produce target text and a special token to in-
 198 dicate the switching from generation to taking
 199 speech input. To this end, we use Llama-3.1-8B-
 200 Instruct (Grattafiori et al., 2024)¹ and employ a
 201 multi-turn dialogue format to formulate the input.
 202 We first feed a system instruction

203 Translate the following speech
 204 from <LangX> to <LangY>.

205 We then add a special USER token to indicate
 206 that the following 12 embeddings and a trailing
 207 END-OF-TURN token are for speech input. We then
 208 prompt the LLM with a special ASSISTANT token
 209 to force LLM to generate tokens. We add a policy
 210 module to check generated tokens. When the pol-
 211 icy module encounters the special END-OF-TURN
 212 token, it will feed a special USER token and take
 213 12 new streaming speech embeddings with a trailing
 214 END-OF-TURN token as new input to the LLM. We
 215 will describe later our inference method to incre-

¹<https://huggingface.co/meta-llama/Llama-3.1-8B-Instruct>

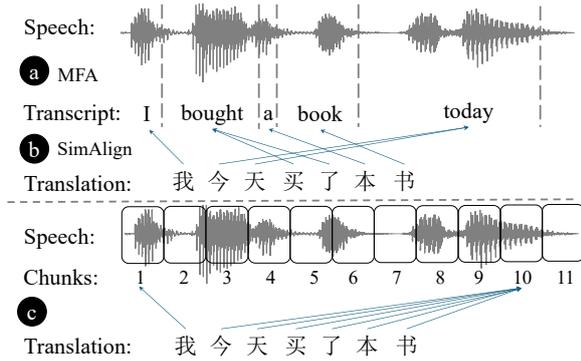


Figure 2: Segmenting speech into chunks and monotonically aligning with translation (bottom).

mentally compute embeddings and generate target tokens for infinite speech input.

3.3 Training Data Construction

SST Trajectory Common speech translation datasets like MuST-C are segmented from complete talks (Di Gangi et al., 2019). To train an SST model in a multi-turn dialogue format, we transform segmented ST triplets (speech s , transcript x , translation y) from MuST-C dataset into SST trajectories. An SST trajectory represents an alternating action sequence of speech reading and translation writing.

As shown in Figure 2, we first align speech utterances with their corresponding transcripts using the Montreal Forced Aligner (MFA) (McAuliffe et al., 2017)². Let m_k^{sx} denote the right boundary of the speech segment corresponding to the transcript token x_k . Also, we utilize SimAlign (Jalili Sabet et al., 2020) with the LaBSE model (Feng et al., 2022) to align words between the transcript and translation. We then monotonize these alignments following Wang et al. (2024). Let $x_{m_i^{xy}}$ be the transcript token that corresponds to translation token y_i . By combining m^{sx} and m^{xy} , we establish a mapping from translation token y_i to its speech boundary $m_i^{sy} = m_{m_i^{xy}}^{sx}$, meaning that y_i is generated after reading $s_{1:m_i^{sy}}$.

Finally, we cut the speech utterance into fixed-length chunks, each lasting 960 ms. We then concatenate translation tokens whose corresponding speech boundaries fall within the same chunk, forming a sequence of trajectory $(s_{C_1}, y_{C_1}), (s_{C_2}, y_{C_2}), \dots$

²<https://github.com/MontrealCorpusTools/Montreal-Forced-Aligner>

Robust Segments for Training Segmented speech utterances primarily consist of human speech; however, non-linguistic sounds (e.g., laughter, applause) are also present. To enhance the robustness of the SST dataset, we cut the entire talk evenly into robust segments that each span 30 speech chunks. If a robust segment starts in the middle of a segmented speech utterance, we shift the robust segment to start with this utterance. The trajectories for a robust segment can then be built by concatenating the trajectories of segmented utterances within this robust segment according to their timestamps and filling the rest translation entries of the trajectory as empty strings.

Multi-Latency Augmentation To further enhance trajectory diversity during training, we propose a simple yet effective multi-latency augmentation strategy. Specifically, given a trajectory $(s_{C_1}, y_{C_1}), (s_{C_2}, y_{C_2}), \dots$, we randomly select a latency multiplier $m \in [1, M]$ and merge every m consecutive chunks of speech with their corresponding translations. The i -th step in the augmented trajectory is then represented as

$$(s_{C_{im}, \dots, C_{(i+1)m-1}}, y_{C_{im}, \dots, C_{(i+1)m-1}}).$$

We also multiply the chunk size of speech encoder with m , i.e., number of frame in a chunk becomes $48m$.

3.4 Training

We train InfiniSST with standard cross-entropy loss on translation tokens, including END-OF-TURN, of the augmented trajectory from robust segments. In the first stage, we freeze the LLM and train only the speech encoder and adapter. In the second stage, we freeze the speech encoder and adapter, training only the LLM.

3.5 Inference on Unbounded Speech

During inference, we cut the unbounded input speech into 960 ms chunks. The latency multiplier m during inference regulates latency by ensuring that translation begins only after every m new chunks have arrived.

At the i -th step, suppose the newly received speech chunks are $C_{im}, \dots, C_{(i+1)m-1}$. Both the speech encoder and the LLM maintain a key-value (KV) cache to prevent redundant computations. Notably, the stored key and value features are extracted *before* applying RoPE, ensuring that no

positional information is embedded within the KV cache.

The speech encoder processes the m new chunks into $48m$ speech features, utilizing the KV cache from chunks $C_{im-w^s+1}, \dots, C_{im-1}$, where w^s is the sliding window size defined in Section 3.2. The adapter then downsamples the $48m$ features into $12m$ embeddings, which are passed to the LLM.

As shown in Figure 3, the LLM employs a sliding window of size w^t . By default, $w^t = 1000$. Inspired by Han et al. (2024), we concatenate the KV cache of instruction with those of the most recent w^t tokens and apply RoPE on top of them. Then the LLM generate translations conditioned on this combined KV cache.

4 Experiment Setups

4.1 Data

We conduct experiments on the En-Es, En-De, and En-Zh directions of the MuST-C dataset (Di Gangi et al., 2019). Due to the poor alignment quality in the En-Zh training set, we filter out misaligned ST triplets using CometKiWi (Rei et al., 2022) and retranslate them using TowerInstruct (Alves et al., 2024). Then, we construct trajectories and robust segments as described in Section 3.3. Further details can be found in the Appendix A.

4.2 Training

We adopt a two-stage supervised fine-tuning approach. In the first stage, we freeze the LLM and train only the speech encoder and adapter for 6 epochs with an effective batch size of 57.6K tokens. We use Adam optimizer (Kingma and Ba, 2017) with learning rate 2×10^{-4} and 1000 warmup steps. We apply gradient clipping with a norm of 1.0. In the second stage, we fine-tune the entire LLM for 1 epoch with an effective batch size of 76.8K tokens and a learning rate of 7×10^{-6} . We employ DeepSpeed Zero Stage-2 optimization³, and enables optimizer and parameter offloading during the second training stage.

4.3 Evaluation

We evaluate SST on complete TED Talks from the MuST-C tst-COMMON set, which consists of 27 TED Talks with durations ranging from 3 to 23 minutes. To assess translation quality, we use SacreBLEU (Post, 2018) and COMET (Guerreiro

et al., 2024). Following the WMT24 practice (Freitag et al., 2024), we compute the COMET score by averaging the scores from XCOMET-XL and XCOMET-XXL. For latency evaluation, we use Length-Adaptive Average Lagging (LAAL) (Papi et al., 2022) for segmented speech baselines and StreamLAAL (Papi et al., 2024a) for unbounded speech, both implemented within the SimulEval framework (Ma et al., 2020a). Computation cost is measured using both computation-aware StreamLAAL (StreamLAAL_CA) and the Real-Time Factor (RTF), defined as the ratio of wall-clock computation time to speech duration.

4.4 Baselines

We compare our method against the following baselines:

AlignAtt (Papi et al., 2023) is a state-of-the-art SST policy applied to offline ST models, translating based on attention scores between translation outputs and speech utterances. It is designed for SST on segmented speech, and we include its results as a reference. We train an offline ST model using segmented ST triplets from MuST-C and robust segments that we constructed. It uses the same model architecture as InfiniSST, except that the speech encoder’s chunk size and window size are set to $+\infty$. We use the attention scores from layer 14 of the LLM and vary the number of frames from 1 to 8.

StreamAtt (Papi et al., 2024a) extends AlignAtt to unbounded speech by maintaining both text and audio history through attention-based selection. We adopt the Fixed-Word approach from StreamAtt, preserving 40 words in the text history. To prevent excessively long preserved speech, we apply truncation when the duration exceeds 28.8 seconds.

StreamAtt+ We observe that the vanilla truncation strategy sometimes removes too much audio, leading to critical misalignment between the preserved speech and its translation. To mitigate this issue, we modify StreamAtt by ensuring that audio segments shorter than 10 seconds are never truncated.

5 Main Results

Competitive Translation Quality at the Same Theoretical Latency Results evaluated with non-computation-aware StreamLAAL are shown in Figure 3. When StreamLAAL is no more than 1.5

³<https://github.com/deepspeedai/DeepSpeed>

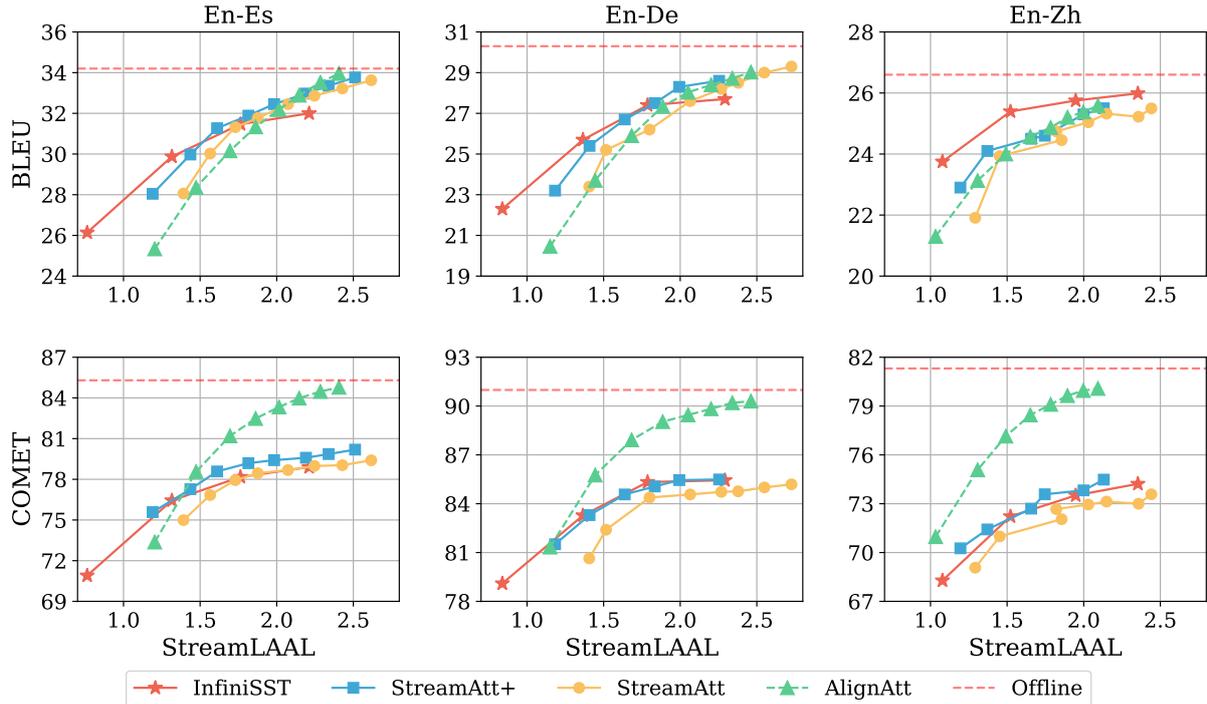


Figure 3: Quality-latency trade-off of InfinitiSST compared to the baselines on complete TED talks from the MuST-C tst-COMMON dataset in the En-Es, En-De, and En-Zh directions. Translation quality is measured using BLEU and COMET scores, while latency is evaluated using the *non-computation-aware* StreamLAAL metric. For reference, we also include offline translation quality and results from AlignAtt tested on segmented speech. InfinitiSST achieves slightly better translation quality than StreamAtt at latency ≤ 1.5 seconds and remains competitive at higher latency levels.

second, InfinitiSST achieves slightly higher BLEU scores (0.5 \sim 1.0) and similar COMET scores than StreamAtt+ on all three language directions. When StreamLAAL is more than 1.5 second, InfinitiSST still achieves higher BLEU score on En-Zh direction and competitive with StreamAtt+ on the En-De and En-Es directions. We note that AlignAtt tested on segmented speech exhibit significant higher COMET scores but not BLEU scores than both InfinitiSST and StreamAtt on all three language directions. A possible reason is that StreamLAAL uses mWERSegmenter (Matusov et al., 2005) to find alignment between translation of the complete talk and segmented references, and COMET is more sensitive to such misalignment than BLEU.

Significantly Lower Computation Cost We run all inference experiments on a single NVIDIA L40S GPU and an AMD EPYC 9354 32-Core CPU. Results evaluated with StreamLAAL_CA are shown in Figure 4. InfinitiSST achieves 0.5 to 1 second lower computation aware latency compared to StreamAtt and StreamAtt+ at the same quality level.

We also compare the Real-Time Factor (RTF) of InfinitiSST and StreamAtt+ in Figure 8. The RTF of InfinitiSST is significantly lower than StreamAtt+, indicating that the computation overhead of InfinitiSST is less than half of the StreamAtt+.

6 Ablation Studies

The default model we use in the ablation study is trained with robust segments and a maximum latency multiplier of $M = 4$ on the En-Zh direction.

6.1 Data

Robust Segments We evaluate the effectiveness of robust segments by comparing InfinitiSST trained on trajectories of robust segments with InfinitiSST trained on trajectories of original MuST-C segmented speech. Both models are evaluated on tst-COMMON En-Zh with latency multipliers $m \in [1, 4]$, and the results are presented in Table 1.

The model trained on trajectories of non-robust segments exhibits abnormal latency scores and lower translation quality compared to the model trained on mega-trajectories. Manual examination of translation instances reveals that the segmented

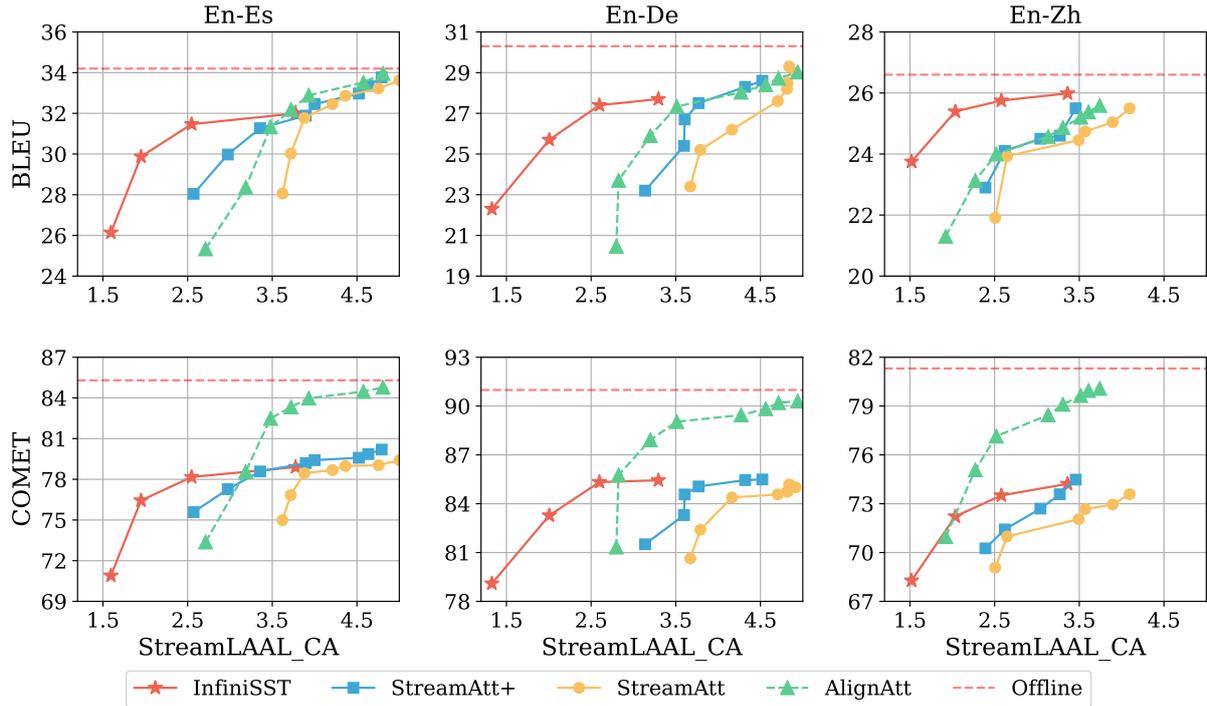


Figure 4: Quality-latency trade-off of InfiniSST compared to the baselines on complete TED talks from the MuST-C tst-COMMON dataset in the En-Es, En-De, and En-Zh directions. Translation quality is measured using BLEU and COMET scores, while latency is evaluated using the *computation-aware* StreamLAAL metric. For reference, we also include offline translation quality and results from AlignAtt tested on segmented speech. InfiniSST achieves significantly lower computation-aware latency compared to StreamAtt at the same quality.

Robust Segments	Non-Robust Segments	Non-Robust Segments*
69.2 / 1.1	50.5 / -220	51.0 / -207
71.9 / 1.5	53.4 / -116	58.1 / -58
72.3 / 1.9	68.4 / 2	65.7 / -22
73.0 / 2.4	66.8 / -12	67.2 / -6

Table 1: Impact of robust segments evaluated on MuST-C En-Zh tst-COMMON with latency multipliers $m = 1, 2, 3, 4$. The model trained on non-robust segments fails to translate unbounded speech. *We suppress the non-linguistic sound tokens but still the model fails to generalize.

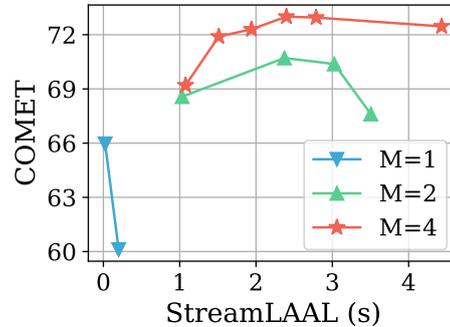


Figure 5: InfiniSST trained with max latency multiplier $M = 1, 2, 4$.

speech model frequently falls into repetition of non-linguistic tokens such as (笑声) whenever non-linguistic sounds appear in the audio.

We attempted to suppress these tokens, and the results are reported in the last column of Table 1. Instead of producing repetitive tokens, the model stops generating translations upon encountering non-linguistic sounds. These findings highlight the importance of training with robust segments.

Multi-Latency We evaluate the effectiveness of the multi-latency augmentation strategy during training. Specifically, we train models with a maximum latency multiplier of $M = 1, 2$, and 4 and perform inference with latency multipliers $m \leq M+2$. The results are presented in Figure 5.

The model trained with $M = 1$ fails to generalize to $m = 2$ during inference, exhibiting a significant drop in translation quality and being unable to generate meaningful translations for larger m . The model trained with $M = 2$ generalizes somewhat

Speech Cache Window w^s	LLM Cache Window w^t	Quality / Latency
10	1000	69.2 / 1.1
5	1000	68.7 / 1.1
20		68.3 / 1.0
40		66.1 / 0.9
10	500	69.0 / 1.0
	2000	69.4 / 1.2
	4000	69.4 / 1.2

Table 2: Impact of cache size during inference. Quality is evaluated with COMET and latency is evaluated with StreamLAAL (unit is second). Model is trained with speech encoder sliding window $w^s = 10$ and no sliding window for LLM. Latency multiplier is set to $m = 1$.

to $m = 3$, but its quality deteriorates significantly when using $m = 4$. In contrast, the model trained with $M = 4$ achieves the best quality-latency trade-off and generalization; however, it still experiences a quality drop when $m = 5$ and 6.

These findings suggest that using a relatively large M during training while ensuring $m \leq M$ during inference is crucial for achieving the best quality-latency trade-off.

6.2 Speech Encoder

Inference Cache Window We first evaluate how the speech encoder’s cache window during inference affects model performance. The model is trained with $w^s = 10$ and tested with $w^s = 5, 10, 20, \text{ and } 40$. The results, presented in Table 2, indicate that using a different cache window size during inference than the one used during training degrades translation quality.

Training Cache Window Furthermore, we train models with different cache window sizes $w^s = 10, 20, 30$ while ensuring that the cache window size matches between training and inference. Since each mega-chunk has a size of 30, training with $w^s = 30$ disables the sliding window mechanism. The results, shown in Figure 7, reveal a surprising observation: the model trained with $w^s = 30$ successfully scales to unbounded speech during inference despite not using a sliding window during training. It also achieves a slightly better quality-latency trade-off compared to the model trained with $w^s = 10$. These findings suggest using the largest possible speech cache window that GPU memory allows.

Model	Talks $\leq 10\text{min}$	Talks $> 10\text{min}$
Llama-3-8K	70.9 / 1.0	67.1 / 1.1
Llama-3.1-128K	71.6 / 1.0	68.0 / 1.1

Table 3: Impact of LLM context length. Llama-3 with 8K context length is still able to generalize to talks longer than 10 minutes.

6.3 LLM

Cache Instruction As described in Section 3.5, we explicitly preserve the KV cache of the translation instruction at the beginning (i.e., the system prompt). If this cache is not retained, the LLM stops translating once the window starts sliding.

Cache Window w^t We evaluate the impact of the LLM’s cache window size during inference on model performance. Notably, the sliding window mechanism is not applied to the LLM during training. We vary the LLM cache window size as $w^t = 500, 1000, 2000, 4000$, and the results are presented in Table 2. Increasing the KV cache size slightly improves translation quality ($69 \rightarrow 69.4$) at the cost of marginally higher latency ($1.0 \rightarrow 1.2$). Compared to the speech encoder, the LLM demonstrates greater robustness to different KV cache window sizes.

Base LLM Context Length Throughout our experiments, we use Llama-3.1-8B-Instruct as the base LLM, which supports a context length of up to 128K tokens. To assess whether InfiniSST generalizes to an LLM with a shorter context limit, we replace it with Llama-3-8B-Instruct, which has an 8K context length⁴. The results, presented in Table 3, indicate that while Llama-3 exhibits lower translation quality compared to Llama-3.1, it is still capable of generalizing to unbounded speech with InfiniSST.

7 Conclusion

We propose InfiniSST that enables simultaneous translation of unbounded speech with state-of-the-art quality latency trade-off on three language directions of MuST-C dataset. Our ablations demonstrate the effectiveness of our carefully constructed data, including robust segments and multi-latency augmentation, and cache management strategy during inference.

⁴A 10-minute speech already generates $10 \cdot 60 \cdot 12 = 7.2K$ speech embeddings, exceeding the 8K context limit of Llama-3-8B-Instruct if combined with the translation tokens.

526 Limitations

527 On the higher theoretical latency level, In-
528 finiSST still falls behind AlignAtt and StreamAtt
529 in some cases. This can be attributed to the lim-
530 ited bidirectional attention of the chunkwise-causal
531 speech encoder. Also, we evaluated on En-X di-
532 rections but not on other directions like X-En and
533 X-X. We have not experimented other pretrained
534 speech encoders and non-Llama LLMs due to com-
535 putation budget. Besides, the StreamLAAL metric
536 is not perfectly reliable due to alignment errors of
537 mWERSegmenter. Finally, we have not conducted
538 human evaluation on user experience of different
539 SST models, which might reveal undetected flaws
540 in current models.

541 References

542 Duarte Miguel Alves, José Pombal, Nuno M Guerreiro,
543 Pedro Henrique Martins, João Alves, Amin Farajian,
544 Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta
545 Agrawal, Pierre Colombo, José G. C. de Souza, and
546 Andre Martins. 2024. [Tower: An open multilingual
547 large language model for translation-related tasks](#). In
548 *First Conference on Language Modeling*.

549 Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed,
550 and Michael Auli. 2020. [wav2vec 2.0: A framework
551 for self-supervised learning of speech representations](#).
552 In *Advances in Neural Information Processing Sys-
553 tems*, volume 33, pages 12449–12460. Curran Asso-
554 ciates, Inc.

555 Keqi Deng, Shinji Watanabe, Jiatong Shi, and Siddhant
556 Arora. 2022. [Blockwise streaming transformer for
557 spoken language understanding and simultaneous
558 speech translation](#). In *Interspeech 2022*, pages 1746–
559 1750.

560 Mattia A. Di Gangi, Roldano Cattoni, Luisa Bentivogli,
561 Matteo Negri, and Marco Turchi. 2019. [MuST-C: a
562 Multilingual Speech Translation Corpus](#). In *Proceed-
563 ings of the 2019 Conference of the North American
564 Chapter of the Association for Computational Lin-
565 guistics: Human Language Technologies, Volume 1
566 (Long and Short Papers)*, pages 2012–2017, Min-
567 neapolis, Minnesota. Association for Computational
568 Linguistics.

569 Domenic Donato, Lei Yu, and Chris Dyer. 2021. [Di-
570 verse pretrained context encodings improve docu-
571 ment translation](#). In *Proceedings of the 59th Annual
572 Meeting of the Association for Computational Lin-
573 guistics and the 11th International Joint Conference
574 on Natural Language Processing (Volume 1: Long
575 Papers)*, pages 1299–1311, Online. Association for
576 Computational Linguistics.

577 Qian Dong, Yaoming Zhu, Mingxuan Wang, and Lei
578 Li. 2022. [Learning when to translate for streaming](#)

[speech](#). In *Proceedings of the 60th Annual Meet-
ing of the Association for Computational Linguistics
(Volume 1: Long Papers)*, pages 680–694, Dublin,
Ireland. Association for Computational Linguistics.

579 Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Ari-
580 vazhagan, and Wei Wang. 2022. [Language-agnostic
581 BERT sentence embedding](#). In *Proceedings of the
582 60th Annual Meeting of the Association for Compu-
583 tational Linguistics (Volume 1: Long Papers)*, pages
584 878–891, Dublin, Ireland. Association for Computa-
585 tional Linguistics.

586 Markus Freitag, Nitika Mathur, Daniel Deutsch, Chi-
587 Kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian
588 Thompson, Frederic Blain, Tom Kocmi, Jiayi Wang,
589 David Ifeoluwa Adelani, Marianna Buchicchio,
590 Chrysoula Zerva, and Alon Lavie. 2024. [Are LLMs
591 breaking MT metrics? results of the WMT24 metrics
592 shared task](#). In *Proceedings of the Ninth Confer-
593 ence on Machine Translation*, pages 47–81, Miami,
594 Florida, USA. Association for Computational Lin-
595 guistics.

596 C. Fugen, M. Kolss, D. Bernreuther, M. Paulik,
597 S. Stuker, S. Vogel, and A. Waibel. 2006. [Open
598 domain speech recognition & translation:lectures and
599 speeches](#). In *2006 IEEE International Conference on
600 Acoustics Speech and Signal Processing Proceedings*,
601 volume 1, pages I–I.

602 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,
603 Abhinav Pandey, Abhishek Kadian, Ahmad Al-
604 Dahle, Aiesha Letman, Akhil Mathur, Alan Schel-
605 ten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh
606 Goyal, Anthony Hartshorn, Aobo Yang, Archi Mi-
607 tra, Archie Sravankumar, Artem Korenev, Arthur
608 Hinsvark, Arun Rao, Aston Zhang, Aurelien Ro-
609 driguez, Austen Gregerson, Ava Spataru, Baptiste
610 Roziere, Bethany Biron, Binh Tang, Bobbie Chern,
611 Charlotte Caucheteux, Chaya Nayak, Chloe Bi,
612 Chris Marra, Chris McConnell, Christian Keller,
613 Christophe Touret, Chunyang Wu, Corinne Wong,
614 Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Al-
615 lonsius, Daniel Song, Danielle Pintz, Danny Livshits,
616 Danny Wyatt, David Esiobu, Dhruv Choudhary,
617 Dhruv Mahajan, Diego Garcia-Olano, Diego Perino,
618 Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy,
619 Elina Lobanova, Emily Dinan, Eric Michael Smith,
620 Filip Radenovic, Francisco Guzmán, Frank Zhang,
621 Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis An-
622 derson, Govind Thattai, Graeme Nail, Gregoire Mi-
623 alon, Guan Pang, Guillem Cucurell, Hailey Nguyen,
624 Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan
625 Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Is-
626 han Misra, Ivan Evtimov, Jack Zhang, Jade Copet,
627 Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park,
628 Jay Mahadeokar, Jeet Shah, Jelmer van der Linde,
629 Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu,
630 Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang,
631 Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park,
632 Joseph Rocca, Joshua Johnstun, Joshua Saxe, Jun-
633 teng Jia, Kalyan Vasuden Alwala, Karthik Prasad,
634 Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth
635 Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer,
636
637
638
639

640	Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal	Elaine Montgomery, Eleonora Presani, Emily Hahn,	704
641	Lakhotia, Lauren Rantala-Yearly, Laurens van der	Emily Wood, Eric-Tuan Le, Erik Brinkman, Este-	705
642	Maaten, Lawrence Chen, Liang Tan, Liz Jenkins,	ban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun,	706
643	Louis Martin, Lovish Madaan, Lubo Malo, Lukas	Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat	707
644	Blecher, Lukas Landzaat, Luke de Oliveira, Madeline	Ozgenel, Francesco Caggioni, Frank Kanayet, Frank	708
645	Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar	Seide, Gabriela Medina Florez, Gabriella Schwarz,	709
646	Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew	Gada Badeer, Georgia Swee, Gil Halpern, Grant	710
647	Oldham, Mathieu Rita, Maya Pavlova, Melanie Kam-	Herman, Grigory Sizov, Guangyi, Zhang, Guna	711
648	badur, Mike Lewis, Min Si, Mitesh Kumar Singh,	Lakshminarayanan, Hakan Inan, Hamid Shojanaz-	712
649	Mona Hassan, Naman Goyal, Narjes Torabi, Niko-	eri, Han Zou, Hannah Wang, Hanwen Zha, Haroun	713
650	lay Bashlykov, Nikolay Bogoychev, Niladri Chatterji,	Habeeb, Harrison Rudolph, Helen Suk, Henry As-	714
651	Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick	pegren, Hunter Goldman, Hongyuan Zhan, Ibrahim	715
652	Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vas-	Phie, Jonathan Torres, Josh Ginsburg, Junjie Wang,	716
653	sic, Peter Weng, Prajjwal Bhargava, Pratik Dubal,	Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khan-	717
654	Praveen Krishnan, Punit Singh Koura, Puxin Xu,	delwal, Katayoun Zand, Kathy Matosich, Kaushik	718
655	Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj	Veeraraghavan, Kelly Michelena, Keqian Li, Ki-	719
656	Ganapathy, Ramon Calderer, Ricardo Silveira Cabral,	ran Jagadeesh, Kun Huang, Kunal Chawla, Kyle	720
657	Robert Stojnic, Roberta Raileanu, Rohan Maheswari,	Huang, Lailin Chen, Lakshya Garg, Lavender A,	721
658	Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ron-	Leandro Silva, Lee Bell, Lei Zhang, Liangpeng	722
659	nie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan	Guo, Licheng Yu, Liron Moshkovich, Luca Wehrst-	723
660	Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sa-	edt, Madian Khabsa, Manav Avalani, Manish Bhatt,	724
661	hana Chennabasappa, Sanjay Singh, Sean Bell, Seo-	Martynas Mankus, Matan Hasson, Matthew Lennie,	725
662	hyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sha-	Matthias Reso, Maxim Groshev, Maxim Naumov,	726
663	ran Narang, Sharath Rapparthi, Sheng Shen, Shengye	Maya Lathi, Meghan Keneally, Miao Liu, Michael L.	727
664	Wan, Shruti Bhosale, Shun Zhang, Simon Van-	Seltzer, Michal Valko, Michelle Restrepo, Mihir Pa-	728
665	denhende, Soumya Batra, Spencer Whitman, Sten	tel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark,	729
666	Sootla, Stephane Collot, Suchin Gururangan, Syd-	Mike Macey, Mike Wang, Miquel Jubert Hermoso,	730
667	ney Borodinsky, Tamar Herman, Tara Fowler, Tarek	Mo Metanat, Mohammad Rastegari, Munish Bansal,	731
668	Sheasha, Thomas Georgiou, Thomas Scialom, Tobias	Nandhini Santhanam, Natascha Parks, Natasha	732
669	Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal	White, Navyata Bawa, Nayan Singhal, Nick Egebo,	733
670	Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh	Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich	734
671	Ramanathan, Viktor Kerkez, Vincent Gouget, Vir-	Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz,	735
672	ginie Do, Vish Vogeti, Vítor Albiero, Vladan Petro-	Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin	736
673	vic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whit-	Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro	737
674	ney Meers, Xavier Martinet, Xiaodong Wang, Xi-	Rittner, Philip Bontrager, Pierre Roux, Piotr	738
675	aofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xin-	Dollar, Polina Zvyagina, Prashant Ratanchandani,	739
676	feng Xie, Xuchao Jia, Xuwei Wang, Yaelle Gold-	Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel	740
677	schlag, Yashesh Gaur, Yasmine Babaei, Yi Wen,	Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu	741
678	Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao,	Nayani, Rahul Mitra, Rangaprabhu Parthasarathy,	742
679	Zacharie Delpierre Coudert, Zheng Yan, Zhengxing	Raymond Li, Rebekkah Hogan, Robin Battey, Rocky	743
680	Chen, Zoe Papakipos, Aaditya Singh, Aayushi Sri-	Wang, Russ Howes, Ruty Rinott, Sachin Mehta,	744
681	vastava, Abha Jain, Adam Kelsey, Adam Shajnfeld,	Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara	745
682	Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand,	Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov,	746
683	Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei	Satadru Pan, Saurabh Mahajan, Saurabh Verma,	747
684	Baevski, Allie Feinstein, Amanda Kallet, Amit San-	Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lind-	748
685	gani, Amos Teo, Anam Yunus, Andrei Lupu, And-	say, Shaun Lindsay, Sheng Feng, Shenghao Lin,	749
686	res Alvarado, Andrew Caples, Andrew Gu, Andrew	Shengxin Cindy Zha, Shishir Patil, Shiva Shankar,	750
687	Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchan-	Shuqiang Zhang, Shuqiang Zhang, Sinong Wang,	751
688	dani, Annie Dong, Annie Franco, Anuj Goyal, Apar-	Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala,	752
689	jita Saraf, Arkabandhu Chowdhury, Ashley Gabriel,	Stephanie Max, Stephen Chen, Steve Kehoe, Steve	753
690	Ashwin Barambe, Assaf Eisenman, Azadeh Yaz-	Satterfield, Sudarshan Govindaprasad, Sumit Gupta,	754
691	dan, Beau James, Ben Maurer, Benjamin Leonhardi,	Summer Deng, Sungmin Cho, Sunny Virk, Suraj	755
692	Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi	Subramanian, Sy Choudhury, Sydney Goldman, Tal	756
693	Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Han-	Remez, Tamar Glaser, Tamara Best, Thilo Koehler,	757
694	cock, Bram Wasti, Brandon Spence, Brani Stojkovic,	Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim	758
695	Brian Gamido, Britt Montalvo, Carl Parker, Carly	Matthews, Timothy Chou, Tzook Shaked, Varun	759
696	Burton, Catalina Mejia, Ce Liu, Changan Wang,	Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai	760
697	Changkyu Kim, Chao Zhou, Chester Hu, Ching-		761
698	Hsiang Chu, Chris Cai, Chris Tindal, Christoph Fe-		762
699	ichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty,		763
700	Daniel Kreymer, Daniel Li, David Adkins, David		764
701	Xu, Davide Testuggine, Delia David, Devi Parikh,		765
702	Diana Liskovich, Didem Foss, Dingkan Wang, Duc		766
703	Le, Dustin Holland, Edward Dowling, Eissa Jamil,		767

768	Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. 2024. The llama 3 herd of models . <i>Preprint</i> , arXiv:2407.21783.	<i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 144–150, Online. Association for Computational Linguistics.	825 826 827 828
780	Nuno M. Guerreiro, Ricardo Rei, Daan van Stigt, Luisa Coheur, Pierre Colombo, and André F. T. Martins. 2024. xcomet: Transparent machine translation evaluation through fine-grained error detection . <i>Transactions of the Association for Computational Linguistics</i> , 12:979–995.	Xutai Ma, Juan Pino, and Philipp Koehn. 2020b. SimulMT to SimulST: Adapting simultaneous text translation to end-to-end simultaneous speech translation . In <i>Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing</i> , pages 582–587, Suzhou, China. Association for Computational Linguistics.	829 830 831 832 833 834 835 836 837
786	Chi Han, Qifan Wang, Hao Peng, Wenhan Xiong, Yu Chen, Heng Ji, and Sinong Wang. 2024. LM-infinite: Zero-shot extreme length generalization for large language models . In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 3991–4008, Mexico City, Mexico. Association for Computational Linguistics.	Evgeny Matusov, Gregor Leusch, Oliver Bender, and Hermann Ney. 2005. Evaluating machine translation output with automatic sentence segmentation . In <i>Proceedings of the Second International Workshop on Spoken Language Translation</i> , Pittsburgh, Pennsylvania, USA.	838 839 840 841 842 843
795	W. Ronny Huang, Shuo yiin Chang, David Rybach, Rohit Prabhavalkar, Tara N. Sainath, Cyril Allauzen, Cal Peysen, and Zhiyun Lu. 2022. E2e segmenter: Joint segmenting and decoding for long-form asr . <i>Preprint</i> , arXiv:2204.10749.	Michael McAuliffe, Michaela Socolof, Sarah Mihuc, Michael Wagner, and Morgan Sonderegger. 2017. Montreal forced aligner: Trainable text-speech alignment using kaldi . In <i>Interspeech 2017</i> , pages 498–502.	844 845 846 847 848
800	Javier Iranzo-Sánchez, Jorge Iranzo-Sánchez, Adrià Giménez, Jorge Civera, and Alfons Juan. 2024. Segmentation-free streaming machine translation . <i>Transactions of the Association for Computational Linguistics</i> , 12:1104–1121.	Siqi Ouyang, Xi Xu, Chinmay Dandekar, and Lei Li. 2024. Fasst: Fast llm-based simultaneous speech translation . <i>Preprint</i> , arXiv:2408.09430.	849 850 851
805	Masoud Jalili Sabet, Philipp Dufter, François Yvon, and Hinrich Schütze. 2020. SimAlign: High quality word alignments without parallel training data using static and contextualized embeddings . In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 1627–1643, Online. Association for Computational Linguistics.	Sara Papi, Marco Gaido, Matteo Negri, and Luisa Bentivogli. 2024a. StreamAtt: Direct streaming speech-to-text translation with attention-based audio history selection . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3692–3707, Bangkok, Thailand. Association for Computational Linguistics.	852 853 854 855 856 857 858 859
812	Diederik P. Kingma and Jimmy Ba. 2017. Adam: A method for stochastic optimization . <i>Preprint</i> , arXiv:1412.6980.	Sara Papi, Marco Gaido, Matteo Negri, and Marco Turchi. 2022. Over-generation cannot be rewarded: Length-adaptive average lagging for simultaneous speech translation . In <i>Proceedings of the Third Workshop on Automatic Simultaneous Translation</i> , pages 12–17, Online. Association for Computational Linguistics.	860 861 862 863 864 865 866
815	Dan Liu, Mengge Du, Xiaoxi Li, Ya Li, and Enhong Chen. 2021. Cross attention augmented transducer networks for simultaneous translation . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 39–55, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	Sara Papi, Peter Polak, Ondřej Bojar, and Dominik Macháček. 2024b. How "real" is your real-time simultaneous speech-to-text translation system? <i>Preprint</i> , arXiv:2412.18495.	867 868 869 870
822	Xutai Ma, Mohammad Javad Dousti, Changhan Wang, Jiatao Gu, and Juan Pino. 2020a. SIMULEVAL: An evaluation toolkit for simultaneous translation . In	Sara Papi, Marco Turchi, and Matteo Negri. 2023. Alignatt: Using attention-based audio-translation alignments as a guide for simultaneous speech translation . In <i>Interspeech 2023</i> , pages 3974–3978.	871 872 873 874
		Matt Post. 2018. A call for clarity in reporting BLEU scores . In <i>Proceedings of the Third Conference on Machine Translation: Research Papers</i> , pages 186–191, Brussels, Belgium. Association for Computational Linguistics.	875 876 877 878 879

880	Ofir Press, Noah A. Smith, and Mike Lewis. 2021. Train short, test long: Attention with linear biases enables input length extrapolation . <i>ArXiv</i> , abs/2108.12409.	936
881		937
882		938
883	Matthew Raffel, Victor Agostinelli, and Lizhong Chen. 2024. Simultaneous masking, not prompting optimization: A paradigm shift in fine-tuning LLMs for simultaneous translation . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 18302–18314, Miami, Florida, USA. Association for Computational Linguistics.	939
884		
885		
886		
887		
888		
889		
890		
891	Ricardo Rei, Marcos Treviso, Nuno M. Guerreiro, Chrysoula Zerva, Ana C Farinha, Christine Maroti, José G. C. de Souza, Taisiya Glushkova, Duarte Alves, Luisa Coheur, Alon Lavie, and André F. T. Martins. 2022. CometKiwi: IST-unbabel 2022 submission for the quality estimation shared task . In <i>Proceedings of the Seventh Conference on Machine Translation (WMT)</i> , pages 634–645, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.	
892		
893		
894		
895		
896		
897		
898		
899		
900		
901	Yi Ren, Jinglin Liu, Xu Tan, Chen Zhang, Tao Qin, Zhou Zhao, and Tie-Yan Liu. 2020. SimulSpeech: End-to-end simultaneous speech to text translation . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 3787–3796, Online. Association for Computational Linguistics.	
902		
903		
904		
905		
906		
907		
908	Felix Schneider and Alexander Waibel. 2020. Towards stream translation: Adaptive computation time for simultaneous machine translation . In <i>Proceedings of the 17th International Conference on Spoken Language Translation</i> , pages 228–236, Online. Association for Computational Linguistics.	
909		
910		
911		
912		
913		
914	Jianlin Su. 2023. Rectified rotary position embeddings . https://github.com/bojone/rerope .	
915		
916	Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. 2024. Roformer: Enhanced transformer with rotary position embedding . <i>Neurocomput.</i> , 568(C).	
917		
918		
919		
920	Jianlin Su, Yu Lu, Shengfeng Pan, Bo Wen, and Yunfeng Liu. 2021. Roformer: Enhanced transformer with rotary position embedding . <i>ArXiv</i> , abs/2104.09864.	
921		
922		
923	Yutao Sun, Li Dong, Barun Patra, Shuming Ma, Shaohan Huang, Alon Benhaim, Vishrav Chaudhary, Xia Song, and Furu Wei. 2023. A length-extrapolatable transformer . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 14590–14604, Toronto, Canada. Association for Computational Linguistics.	
924		
925		
926		
927		
928		
929		
930		
931	Minghan Wang, Thuy-Trang Vu, Yuxia Wang, Ehsan Shareghi, and Gholamreza Haffari. 2024. Conversational simulmt: Efficient simultaneous translation with large language models . <i>Preprint</i> , arXiv:2402.10552.	
932		
933		
934		
935		
	Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. 2024. Efficient streaming language models with attention sinks . <i>Preprint</i> , arXiv:2309.17453.	940
		941
		942
		943
		944
		945
		946
		947
	Xi Xu, Siqi Ouyang, Brian Yan, Patrick Fernandes, William Chen, Lei Li, Graham Neubig, and Shinji Watanabe. 2024. CMU’s IWSLT 2024 simultaneous speech translation system . In <i>Proceedings of the 21st International Conference on Spoken Language Translation (IWSLT 2024)</i> , pages 154–159, Bangkok, Thailand (in-person and online). Association for Computational Linguistics.	948
		949
		950
		951
	Takeonori Yoshimura, Tomoki Hayashi, Kazuya Takeda, and Shinji Watanabe. 2020. End-to-end automatic speech recognition integrated with ctc-based voice activity detection . <i>Preprint</i> , arXiv:2002.00551.	952
		953
		954
		955
		956
	Donglei Yu, Yang Zhao, Jie Zhu, Yangyifan Xu, Yu Zhou, and Chengqing Zong. 2025. SimulPL: Aligning human preferences in simultaneous machine translation . In <i>The Thirteenth International Conference on Learning Representations</i> .	957
		958
		959
		960
		961
		962
	Xingshan Zeng, Liangyou Li, and Qun Liu. 2021. Real-Trans: End-to-end simultaneous speech translation with convolutional weighted-shrinking transformer . In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 2461–2474, Online. Association for Computational Linguistics.	

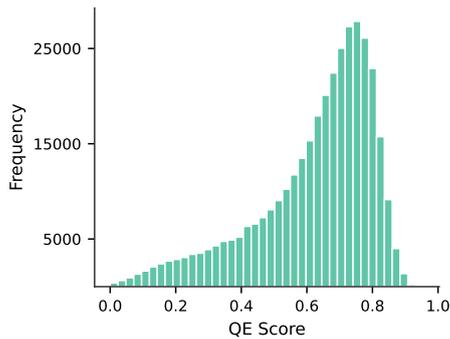


Figure 6: COMET-KIWI quality estimation score distribution on MuST-C en-zh data

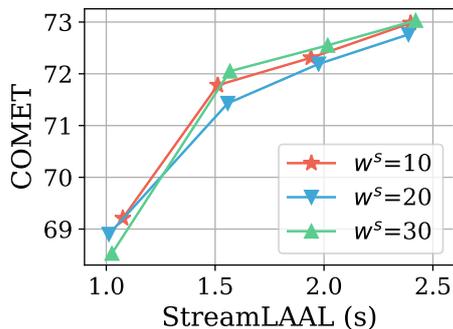


Figure 7: Impact of training-time sliding window size w^s of speech encoder.

A Additional Data Details

A.1 QE filtering and forward translation

We first use Whisper to perform automatic speech recognition (ASR) on all training segments. We then apply CometKiwi⁵ to estimate the quality of ASR outputs by computing quality estimation (QE) scores between the ASR results and the reference text. As shown in Figure 6, we retain only instances where the QE score is greater than 0.5, which accounts for 78.64% of the data, resulting in a total of 280K instances.

Upon further inspection, we observed that many filtered-out cases exhibited acceptable word error rates (WER) between the ASR outputs and the source text. To recover these cases, we performed forward translation using the 7B version of TowerInstruct⁶ on the source text using TowerInstruct with the following decoding settings: temperature = 0.0 and frequency penalty = 0.1. The translations were generated using vLLM.

⁵<https://huggingface.co/Unbabel/wmt23-cometkiwi-da-xxl>

⁶<https://huggingface.co/Unbabel/TowerInstruct-7B-v0.2>

A.2 Dataset Statistics

The MuST-C dataset used in our experiments consists of 105,647 instances for En-Zh, 88,725 for En-Es, and 70,037 for En-De.

Figure 9 shows the reference length distribution across these language pairs.

For En-Zh, the reference text length averages 124.32 characters, with a maximum of 444. En-Es has significantly longer references, averaging 400.06 characters and reaching a maximum of 1,116. En-De also exhibits long references, with an average of 419.47 characters and a maximum of 957. Spanish reference lengths in word count average 67.1 words, with a median of 70.0 and a 90th percentile of 90.0. German references are slightly shorter, averaging 63.6 words, with a median of 65.0 and a 90th percentile of 87.0.

En-Zh segments average 26.85 seconds, En-Es 25.25 seconds, and En-De 26.11 seconds, all with a maximum of 28.80 seconds, reflecting speech-text alignment across languages.

B Additional Experiment Results

Impact of speech encoder window size during training w^s is shown in Figure 7. The RTF of InfiniS-Stand baseline StreamAtt+ is shown in Figure 8.

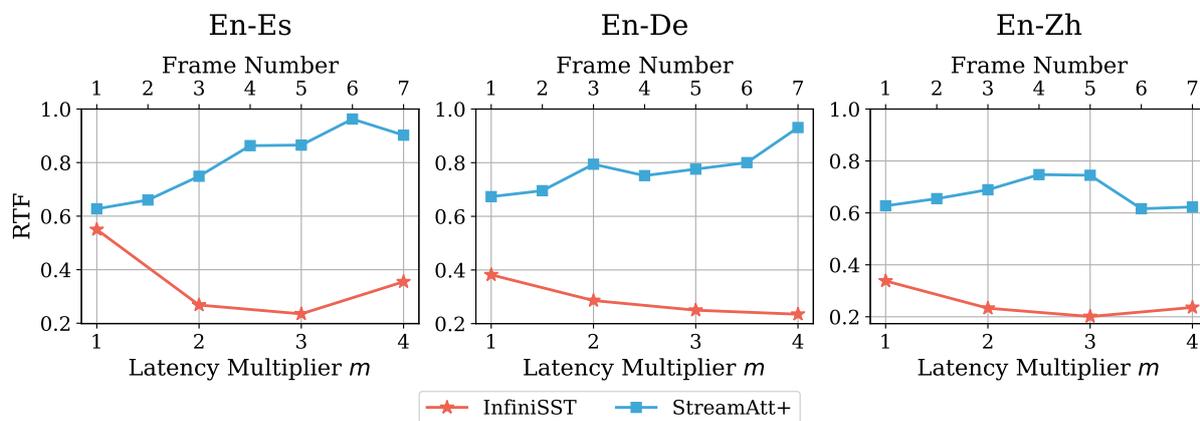


Figure 8: The Real-Time-Factor of InfiniSST and baseline StreamAtt+.

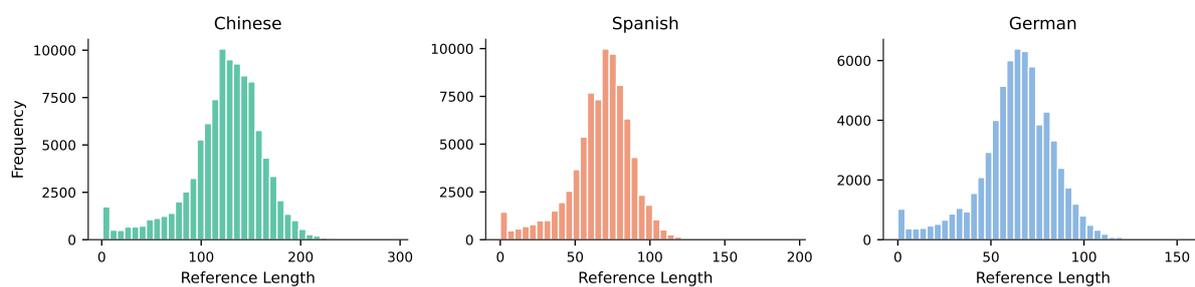


Figure 9: Reference length distribution of SST trajectories on MuST-C En-Zh, En-Es, and En-De.