

# StreamUni: Achieving Streaming Speech Translation with a Unified Large Speech-Language Model

Anonymous ACL submission

## Abstract

Streaming speech translation (StreamST) requires determining appropriate timing, known as policy, to generate translations while continuously receiving source speech inputs, balancing low latency with high translation quality. Existing StreamST methods typically operate on sentence-level speech segments, referred to as simultaneous speech translation (SimulST). In practice, they require collaboration with upstream segmentation models to accomplish StreamST, where the truncated speech segments constrain SimulST models to make policy decisions and generate translations based on pre-defined contextual information preset by the upstream models. Moreover, SimulST models struggle to learn effective policies due to the complexity of speech inputs and cross-lingual generation. To address these challenges, we propose StreamUni, which achieves StreamST through a unified Large Speech-Language Model (LSLM). Specifically, StreamUni incorporates speech Chain-of-Thought (CoT) in guiding the LSLM to generate multi-stage outputs. Leveraging these multi-stage outputs, StreamUni simultaneously accomplishes speech segmentation, policy decision, and translation generation. Additionally, we propose a streaming CoT training method that enhances low-latency policy decisions and generation capabilities using limited CoT data. Experiments demonstrate that our approach achieves state-of-the-art performance on both SimulST and StreamST tasks.

## 1 Introduction

Streaming speech translation (StreamST) (Ma et al., 2019, 2020b; Dong et al., 2022), known as simultaneous interpretation, generates translations while continuously receiving speech inputs. It is commonly employed in cross-lingual scenarios such as international conferences and real-time subtitles.

Unlike traditional offline speech translation (Gangi et al., 2019; Alinejad and Sarkar, 2020;

Lee et al., 2022), StreamST must ensure translation quality and minimize the latency between receiving speech inputs and generating translations (Zhang et al., 2024a). To this end, StreamST requires a generation policy to determine the appropriate timing for outputting translated words. Considering that StreamST is often deployed in scenarios lasting tens of minutes to several hours (Ma et al., 2019), and that the relevant content attended to by StreamST is primarily concentrated around real-time inputs (Papi et al., 2024), it is necessary to implement a truncation policy that can truncate historical speech inputs and translations, enabling the model to focus on recent speech inputs while preventing information overload that could compromise efficiency. Therefore, an ideal StreamST model requires both an effective generation policy and truncation policy to achieve low latency and high translation quality.

Existing methods primarily focus on Simultaneous Speech Translation (SimulST) rather than StreamST. They are typically restricted to short speech clips segmented by Voice Activity Detection (VAD) (Team, 2024; Tang et al., 2023), making them unsuitable for long-form streams. Within SimulST, research centers on generation policies: fixed policies (Ma et al., 2019, 2020b) alternate between fixed-duration chunks and predetermined outputs, often incurring high latency or poor quality by ignoring content; in contrast, adaptive policies—utilizing integrate-and-fire (Dong et al., 2022), CTC (Zhang et al., 2024a), or Transducers (Tang et al., 2023)—dynamically adjust to text density for better performance. Nevertheless, these methods still yield suboptimal quality due to their reliance on small-scale Transformer architectures (Vaswani et al., 2017).

More recent work attempts to leverage Large Speech-Language Models (LSLMs) for SimulST, delivering superior performance. These methods adopt either fixed (Agostinelli et al., 2024) or adap-

085 tive policies, the latter typically requiring fine-  
086 tuning on specialized policy-specific data for au-  
087 toregressive prediction (Wang et al., 2024; Cheng  
088 et al., 2024; Labiausse et al., 2025). However,  
089 such fine-tuning methods compromise the inher-  
090 ent generation capabilities of LSLMs and present  
091 difficulties in efficiently transferring to newly ad-  
092 vanced LSLMs. Furthermore, current research has  
093 inadequately explored truncation policies, with at-  
094 tempts to truncate historical translations through  
095 constructing complex translation trajectory training  
096 data and sliding window schemes (Ouyang et al.,  
097 2025). This approach incurs substantial data con-  
098 struction costs. Consequently, a unified LSLM  
099 paradigm for efficient StreamST has emerged as a  
100 promising alternative.

101 Despite its advantages, implementing StreamST  
102 using a unified LSLM remains challenging, as it  
103 requires LSLM to simultaneously handle trunca-  
104 tion and generation policies while achieving real-  
105 time translation. To determine generation policy,  
106 LSLMs need to detect valid content in real-time  
107 speech stream and decide on the optimal gener-  
108 ation timing and output translations (Dong et al.,  
109 2022). As the speech stream grows, LSLMs require  
110 the truncation policy to discard historical speech  
111 segments and translations, ensuring the model fo-  
112 cuses on recent inputs while avoiding excessive  
113 computational overhead (Papi et al., 2024). Trunca-  
114 tion policy must ensure that discarded speech seg-  
115 ment is fully translated and that discarded transla-  
116 tions accurately correspond to the discarded speech  
117 segments, thereby maintaining truncation integrity.  
118 Beyond policy decisions, StreamST also needs to  
119 accomplish high-quality translation for continu-  
120 ously incoming speech input streams. However,  
121 conventional approaches that separately optimize  
122 these three subtasks require constructing substan-  
123 tial amounts of corresponding training data (Wang  
124 et al., 2024), which is not only resource-intensive  
125 but also present significant difficulties in transfer-  
126 ring to newly advanced LSLMs. Therefore, investi-  
127 gating how to enable LSLMs to efficiently accom-  
128 plish all subtasks in a unified manner for effective  
129 StreamST is of paramount importance.

130 To address these challenges, we propose Stream-  
131 Uni, a framework that efficiently enables a uni-  
132 fied LSLM to accomplish all subtasks of StreamST  
133 in a cohesive manner. StreamUni introduces the  
134 speech Chain-of-Thought (CoT) (Huang et al.,  
135 2023; Nguyen et al., 2024) that guides LSLMs to  
136 progressively generate transcriptions and transla-

137 tions based on the speech inputs. Leveraging multi-  
138 stage outputs, the model handles generation policy,  
139 truncation policy, and streaming translation genera-  
140 tion subtasks. For the generation policy, StreamUni  
141 detects effective speech chunks in real-time to de-  
142 termine optimal generation timing based on the  
143 coherence between the transcription and previously  
144 output translations. For truncation policy, Stream-  
145 Uni maintains transcription queues across differ-  
146 ent timestamps and determines speech truncation  
147 timing by comparing current and historical trans-  
148 criptions. When the generation policy is triggered,  
149 the model generates translations in real-time to en-  
150 sure a seamless interactive experience. When the  
151 truncation policy is activated, the model promptly  
152 truncates historical speech data and correspond-  
153 ing translations in real-time, thereby maintaining  
154 high efficiency. Through this design, StreamUni  
155 achieves StreamST via multi-task results across  
156 multiple stages of the speech CoT.

157 To further enhance streaming performance, we  
158 propose a Streaming CoT training scheme that op-  
159 timizes multi-stage CoT outputs. Therefore, Stream-  
160 Uni unifies all subtasks through the speech CoT  
161 and achieves holistic optimization via a unified  
162 training strategy. Experiments demonstrate that our  
163 method efficiently achieves state-of-the-art perfor-  
164 mance on both SimulST and StreamST tasks.

## 165 2 Background

166 **Streaming Speech Translation** Let the complete  
167 speech stream be represented as  $\mathbf{s} = (s_1, \dots, s_N)$ ,  
168 where  $s_i$  denotes a speech chunk of predefined size,  
169 typically around 640ms. Given the continuously  
170 arriving speech chunks, StreamST progressively  
171 generates translation  $\mathbf{y} = (y_1, \dots, y_I)$  under a gen-  
172 eration policy  $\mathbf{g} = (g_1, \dots, g_I)$  where  $g_i$  represents  
173 the number of speech chunks received when gener-  
174 ating  $y_i$ . Thus, StreamST can be formulated as:

$$175 p(\mathbf{y} | \mathbf{s}, \mathbf{g}) = \prod_{i=1}^I p(y_i | \mathbf{s}_{\leq g_i}, \mathbf{y}_{< i}). \quad (1)$$

176 However, when the incoming speech stream be-  
177 comes excessively long, StreamST models need to  
178 truncate historical speech and translations in real-  
179 time, thereby focusing on recent inputs while avoid-  
180 ing excessive inference latency (Iranzo-Sánchez  
181 et al., 2024). Consequently, truncation policy is  
182 employed to determine truncation timing. Let  
183 the truncation policy for the overall speech input  
184 and target translation be  $\mathbf{a} = (a_1, \dots, a_M)$  and

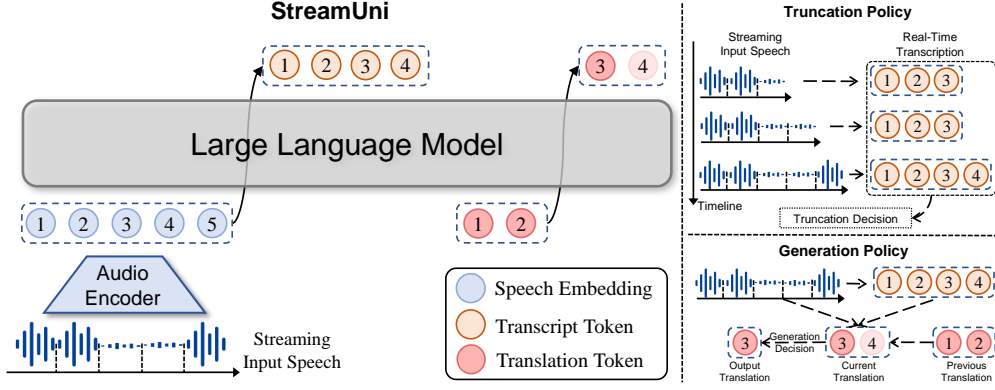


Figure 1: The framework of StreamUni and illustration of truncation policy and generation policy. The model will generate real-time transcription based on the existing speech inputs and compare it with historical transcriptions to determine the truncation policy. If truncation is decided, the model will bypass the generation policy and continue generating the full translation of the current speech. Otherwise, the model will determine the number of translation words to continue generating based on the lag relationship between the real-time transcription and the output translation, and generate the translation using CoT.

$\mathbf{b} = (b_1, \dots, b_M)$  respectively, where  $M$  denotes the desired number of truncated segments, and  $a_m$  and  $b_m$  represent the ending positions of the  $m$ -th segment within the complete input stream and translation. Under the guidance of the segmentation policy, StreamST is reformulated as:

$$p(\mathbf{y} \mid \mathbf{s}, \mathbf{g}, \mathbf{a}, \mathbf{b}) = \prod_{i=1}^{b_1} p(y_i \mid \mathbf{s}_{1:g_i}, \mathbf{y}_{1:i-1}) \times \prod_{m=2}^M \prod_{i=b_{m-1}+1}^{b_m} p(y_i \mid \mathbf{s}_{a_{m-1}+1:g_i}, \mathbf{y}_{b_{m-1}+1:i-1}) \quad (2)$$

where the streaming translation generation will be based solely on the input speech segment and output translation segment that remain after truncation. Therefore, StreamST requires determining both truncation and generation policies to guide the model in accomplishing translation generation.

**Chain-of-Thought** Chain-of-Thought (CoT) is originally developed for text-based tasks and has been proven to enhance performance on complex tasks by prompting large language models (LLMs) to think step by step before giving final results (Wei et al., 2022; DeepSeek-AI et al., 2025a). For speech inputs, CoT techniques have been widely adopted in speech-to-text cross-modal tasks, where LSLMs first generate transcription and subsequently produce the final outputs (Zhang et al., 2023; Huang et al., 2023). In the context of speech translation, the model first generates transcription  $\mathbf{x} =$

$(x_1, \dots, x_J)$ , followed by the translation:

$$p(\mathbf{y} \mid \mathbf{s}) = p(\mathbf{y} \mid \mathbf{x}, \mathbf{s})p(\mathbf{x} \mid \mathbf{s}). \quad (3)$$

### 3 Method

In this section, we propose StreamUni, a framework that leverages speech CoT to consolidate all subtasks in StreamST. We begin by introducing the architecture of StreamUni and detailing its operational process for achieving StreamST. Subsequently, we present the truncation and generation policies within the StreamUni framework, which governs the management of historical speech and translation and the real-time translation generation. To further enhance the generation capabilities of LSLMs across multiple CoT stages under low-latency conditions, we propose a novel streaming CoT training scheme. The following subsections detail our methodology.

#### 3.1 Model Framework

The framework is illustrated in Figure 1. StreamUni first transcribes the input speech and compares the real-time transcription with the historical transcriptions to determine the truncation policy. If the truncation policy is triggered, StreamUni directly generates the translation bypassing the generation policy; otherwise, the number of words to be generated is determined by the generation policy and translation is generated accordingly. A formalized operational process is presented as follows.

Given that the previous truncation timing of the input speech stream is  $a_m$ , and the current timing is  $n$  ( $n > a_m$ ), the currently received speech

segment fed into the model can be represented as  $\mathbf{s}_{a_m+1:n}$ . For segment  $\mathbf{s}_{a_m+1:n}$ , the LSLM first utilizes an audio encoder to encode it into speech embeddings. Following the speech CoT instruction, LSLM subsequently generates real-time transcription  $\mathbf{x}^{(n)}$  of  $\mathbf{s}_{a_m+1:n}$ . StreamUni then determines the truncation policy by comparing  $\mathbf{x}^{(n)}$  with maintained historical transcription queue, specifically deciding whether the current timing  $n$  should trigger truncation. If it is determined that the current timing  $n$  should trigger truncation, StreamUni disregards the generation policy, and continues generating and outputting all subsequent translation based on the input segment  $\mathbf{s}_{a_m+1:n}$  and real-time transcription  $\mathbf{x}^{(n)}$ , building upon the already output translation segment  $\mathbf{y}_{b_m+1:i-1}$ , where  $b_m$  is the translation truncation index corresponding to  $a_m$ . Otherwise, StreamUni determines the generation policy based on the real-time transcript  $\mathbf{x}^{(n)}$  and uses it to determine the number of output words at current timing. We then elaborate on the truncation policy and generation policy in detail.

**Truncation Policy** StreamST employs a truncation policy to remove historical speech and translation segments no longer required for subsequent generation. To ensure truncation integrity, each truncated speech segment must maintain semantic alignment with its corresponding translation segment (Iranzo-Sánchez et al., 2024). The above truncation constraints serve dual purposes: (1) preventing the eliminated speech segment containing untranslated content, which would compromise generation quality, and (2) avoiding removal of already-translated content of remaining speech segment, which will result in repetitive translation of remaining speech segment. According to these, we propose the following truncation policy.

For speech stream  $\mathbf{s} = (s_1, \dots, s_N)$ , StreamUni obtains real-time transcription after receiving each chunk and maintains a historical transcription queue  $\mathbf{q}$ . Assuming the end position of the previous truncated input segment is  $a_m$  and the chunk to be processed is  $n$  ( $n > a_m$ ),  $\mathbf{q}$  can be represented as  $[\mathbf{x}^{(a_m+1)}, \dots, \mathbf{x}^{(n-1)}]$ . StreamUni first obtains the transcription  $\mathbf{x}^{(n)}$  based on  $\mathbf{s}_{a_m+1:n}$ :

$$\mathbf{x}^{(n)} = \arg \max_{\mathbf{x}} p(\mathbf{x} | \mathbf{s}_{a_m+1:n}). \quad (4)$$

Subsequently, we compare  $\mathbf{x}^{(n)}$  with items in  $\mathbf{q}$  to determine the truncation policy. Speech segment truncation occurs if either condition is satisfied:

- If  $\mathbf{x}^{(n)}$  remains identical to real-time transcriptions from the previous two chunks ( $\mathbf{x}^{(n-1)}$  and  $\mathbf{x}^{(n-2)}$ ), then  $a_{m+1} = n$  becomes the speech truncation timing and  $\mathbf{s}_{a_m+1:a_{m+1}}$  is discarded. The historical transcription queue is cleared ( $\mathbf{q} = []$ ).
- If  $\mathbf{x}^{(l)}$  ( $l = n-1, n-2$ ) forms a complete sentence terminated by punctuation (.?!), and  $\mathbf{x}^{(n)}$  begins a new sentence following the complete sentence, then truncation timing is  $a_{m+1} = l$  and  $\mathbf{s}_{a_m+1:a_{m+1}}$  is discarded. The historical transcription queue is cleared, and the newly generated  $\mathbf{x}^{(l)}$  ( $l = a_{m+1} + 1, \dots, n$ ) are sequentially added.

After determining the truncation timing, StreamUni generates and outputs the complete translation corresponding to  $\mathbf{s}_{a_m+1,a_{m+1}}$  based on previously output translation:

$$\mathbf{y}_{i:b_{m+1}} = \arg \max_{\mathbf{y}} p(\mathbf{y} | \mathbf{s}_{a_m+1:a_{m+1}}, \mathbf{x}^{(a_m+1)}, \mathbf{y}_{b_m+1:i-1}), \quad (5)$$

where  $b_{m+1}$  is the index of the last word in the output translation. The translation segment  $\mathbf{y}_{b_m+1:b_{m+1}}$  is then discarded.

In conclusion, we select truncation timing when users maintain prolonged silence or finish a full sentence, as content prior to this timing is relatively complete and subsequent translations are unlikely to reference earlier inputs. After determining input truncation timing, target truncation timing is decided by outputting complete translation for the truncated input segment, thereby maintaining semantic integrity of the truncation. More explanation is in Appendix A.

**Generation Policy** After establishing the truncation policy, we then determine the generation policy, which controls model output at all timing except truncation moments. The generation policy follows two key principles. First, the model should continue generating translation upon detecting the text within input speech; otherwise, no generation is required (Dong et al., 2022). Second, translation generation should lag behind the valid source inputs to provide sufficient context for translation (Liu et al., 2021). Leveraging speech CoT, we implement the generation policy in Figure 1.

Assume the previous truncated segment is the  $m$ -th segment, and the speech chunk to be processed

is  $s_n$ . We can obtain the transcription  $\mathbf{x}^{(n)}$  using Eq.(4). Let  $C$  denote the number of words in  $\mathbf{x}^{(n)}$  and  $i-1$  represent the position of the last word in the already output translation. The number of translation words allowed to be output is:

$$O = C - k - (i - 1 - b_m), \quad (6)$$

where the second term  $k$  is the delay hyperparameter, and the third term represents the number of retained output translation words. This setting ensures that translation generation consistently lags behind the input text by  $k$  words, providing sufficient context for generation. The current translation generation can be represented as:

$$\mathbf{y}_{i:i-1+O} = \arg \max_{\mathbf{y}} p(\mathbf{y} | \mathbf{s}_{a_m+1:n}, \mathbf{x}^{(n)}, \mathbf{y}_{b_m+1:i-1}). \quad (7)$$

Then the translation  $\mathbf{y}_{i:i-1+O}$  will be output.

### 3.2 Streaming CoT Training

After introducing the overall framework, StreamUni can now perform StreamST using existing LSLMs (Microsoft et al., 2025; Xu et al., 2025). However, existing LSLMs are trained on multi-task datasets containing complete speech inputs paired with corresponding responses. In streaming scenarios with continuously growing speech stream, LSLMs must handle speech inputs of different lengths, which we refer to as streaming generation capability. Furthermore, our approach unifies policy decisions and streaming translation generation through speech CoT, which requires enhanced streaming generation capability across multiple stages of speech CoT. Therefore, we propose the Streaming CoT training scheme, which improves the capabilities of policy decision and streaming translation generation by augmenting streaming speech CoT data.

Our method constructs streaming CoT data using existing non-streaming CoT triplets of speech, transcription, and translation. Given the input speech stream  $\mathbf{s} = (s_1, \dots, s_N)$ , our approach randomly truncates the stream through uniform sampling to obtain  $\mathbf{s}_{\leq i}$ . We then employ timestamp alignment tools to extract the corresponding transcription  $\mathbf{x}^{(i)}$  for  $\mathbf{s}_{\leq i}$  from the complete transcription  $\mathbf{x}$ . Our Streaming CoT training encourages the LSLM to predict full translation based on partial speech and

transcription:

$$\mathcal{L} = - \sum_{\mathbf{s}_{\leq i} \sim \mathcal{U}(\mathbf{S})} \log p(\mathbf{y} | \mathbf{x}^{(i)}, \mathbf{s}_{\leq i}) p(\mathbf{x}^{(i)} | \mathbf{s}_{\leq i}), \quad (8)$$

where  $\mathbf{S}$  is  $\{\mathbf{s}_{\leq 1}, \dots, \mathbf{s}_{\leq N}\}$  and  $\mathbf{s}_{\leq i} \sim \mathcal{U}(\mathbf{S})$  represents uniform sampling from set  $\mathbf{S}$ . This formulation trains accurate transcription prediction for policy decisions while requiring complete translation prediction to enhance generation capability and prevent premature termination. For efficiency, we employ sampling rather than training on all possible speech inputs for an instance. Through this training approach, our method efficiently enhance streaming CoT generation capability, thereby improving the capabilities of policy decision and streaming translation generation in low latency. In experiments, our training method requires integration with traditional non-streaming training approaches to achieve greater performance gains.

## 4 Experiments

### 4.1 Datasets

We mainly conduct experiments on streaming speech translation (StreamST) and simultaneous machine translation (SimulST) tasks.

**MuST-C English  $\Rightarrow$  German (En  $\Rightarrow$  De)** This dataset (Di Gangi et al., 2019) is collected from TED talks. The dataset contains both document-level and sentence-level speech translation data, enabling evaluation of SimulST and StreamST tasks.

**MuST-C English  $\Rightarrow$  Spanish (En  $\Rightarrow$  Es)** The dataset is constructed following the same approach as MuST-C En-De and serves as an evaluation benchmark for both StreamST and SimulST tasks.

**CoVoST2 English  $\Rightarrow$  Chinese (En  $\Rightarrow$  Zh)** This dataset only contains sentence-level speech translation data and is used to evaluate SimulST tasks (Wang et al., 2020).

**CoVoST2 French  $\Rightarrow$  English (Fr  $\Rightarrow$  En)** This dataset is also used to evaluate SimulST tasks.

### 4.2 System Settings

In this subsection, we delineate the settings of our StreamUni method and then present the comparative methods for each task separately.

For our approach, we adopt Phi-4-Multimodal (Microsoft et al., 2025) as the primary backbone LSLM and fine-tune it using the speech CoT data across four language directions. Specifically, the En  $\Rightarrow$  Zh direction contains 50 hours of streaming

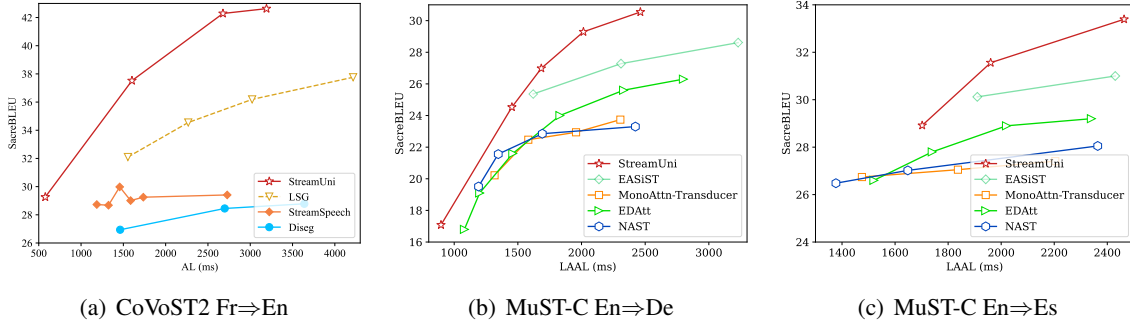


Figure 2: Performance of different methods on SimulST task.

CoT data and 50 hours of non-streaming CoT data, while the other three directions each comprise 100 hours of non-streaming CoT data. The CoT instruction used for LSLM inference is: "Transcribe the audio to text, and then translate the audio to **{target\_lang}**. Use <sep> as a separator between the original transcript and the translation". During inference, the chunk size is set to 320ms for the En-Zh direction and 640ms for the other directions. To control inference latency, we configure  $k$  as  $\{1, 3, 5, 7, 9\}$ . When applied to the SimulST task, StreamUni executes only the generation policy. **Additional training hyperparameters are provided in the Appendix B.** Beyond Phi-4-Multimodal, we also experiment with Qwen2.5-Omni (Xu et al., 2025) as the base LSLM to validate the generalizability of our method, leveraging its thinker for policy-decision and translation generation.

For SimulST task, we compare our method with **DiSeg** (Zhang and Feng, 2023), **NAST** (Ma et al., 2023), **EDAtt** (Papi et al., 2023a), **StreamSpeech** (Zhang et al., 2024b), **LSG** (Guo et al., 2025), **MonoAttn-Transducer** (Ma et al., 2025) and **EASiST** (Fu et al., 2025). We also design a baseline called **Phi4-Wait- $k$** , which uses fine-tuned Phi-4-MM as our StreamUni but employs a generation policy that waits for  $k-1$  chunks and then outputs one word for each subsequently received chunk.

For the StreamST task, we compare our method with **StreamAttFW** and **StreamAttP** (Papi et al., 2024). Furthermore, we implement a baseline called **Phi-4-VAD**, which replaces our truncation policy with VAD (Team, 2024) while keeping all other components consistent with our approach. We also propose an additional cascaded method named **Phi-4-Cascaded**: we adopt Whisper-Large-V3 (Radford et al., 2022) as the ASR model and feed its outputs into Phi-4-Mini-Instruct for translation. The prompt used for Phi-4-Mini-Instruct is:

"Translate the English text to German based on the given German translation. English text: {cot\_asr}. German translation: {cot\_st}."

### 4.3 Evaluation

In evaluating streaming generation systems, we need to assess two critical aspects: latency and generation quality. To quantify latency, we utilize the Average Lagging (AL) (Ma et al., 2019) and Length-Adaptive Average Lagging (LAAL) metrics (Papi et al., 2022a), which measures the delay between input reception and output generation. For translation quality, we use the SacreBLEU (Post, 2018) and COMET (Rei et al., 2022) metrics. For the SimulST task, we employ the SimulEval tool (Ma et al., 2020a) to evaluate our StreamUni. In the StreamST task, we follow the setup of Papi et al. (2024). We first use mWERSegmenter (Matusov et al., 2005) for aligning document-level translation with references and then convert these alignments into consistent metrics used in the SimulST task. In this task, The latency metric is termed StreamLAAL, and translation quality is evaluated using Stream SacreBLEU by comparing the segmented document-level translations with the reference translations.

### 4.4 Main Results

We evaluate our methods on both SimulST and StreamST tasks.

As illustrated in Figure 2, our method achieves optimal SimulST performance across all datasets. Compared to traditional SimulST methods employing Encoder-Decoder architectures (e.g., NAST and EDAtt), our method harnesses the comprehension and generation capabilities of LSLMs (Microsoft et al., 2025), yielding substantial improvements across all latency settings. Although methods like LSG also leverage LSLMs and demon-

strate promising results, their policy decisions rely on heuristic rules (Guo et al., 2025), resulting in suboptimal performance. Furthermore, compared to EASiST built on larger-scale foundation models as backbones and trained with customized policy data (Fu et al., 2025), our model achieves equally better performance. This is mainly attributed to our method’s utilization of intermediate outputs from speech CoT, which enables real-time detection of valid user inputs and allows generation decisions to be made at optimal timings. Through our superior policy and streaming CoT training scheme, we achieve further performance gains.

Our method also demonstrates superior performance on the StreamST task, as shown in Figure 3. Traditional StreamST approaches, including StreamFW and StreamAttP (Papi et al., 2024), rely on attention interpretability to determine generation and truncation policies. In contrast, our approach utilizes speech CoT for real-time detection of valid speech inputs to inform generation policies, while implementing truncation policies through alignments between speech input and translations. This design enables more effective policy decisions and enhanced performance. Compared to Phi-4-VAD, which employs VAD for truncation policy, our method achieves truncation policy through the semantic alignments between speech inputs and translation, resulting in more appropriate timing and enhanced performance. Our method also significantly outperforms Phi-4-Cascaded, particularly in latency, underscoring the superiority of end-to-end models. We anticipate further gains for Phi-4-Mini-Instruct with dedicated translation training.

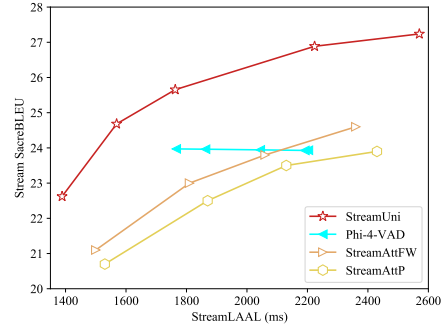
In addition to the above results, we also considered the usability of our method. To this end, we additionally incorporated the actual machine inference time into the latency metric when calculating latency, which is denoted as Computation-Aware LAAL (Xu et al., 2024). It essentially measures the average delay from the user’s speech input to the machine’s output of the corresponding translation. For details, refer to Appendix C.

## 5 Analysis

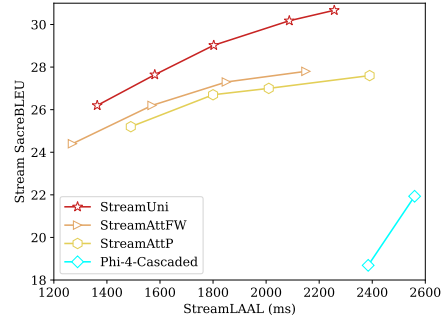
To provide deeper understandings into our approach, we conduct comprehensive analyses.

### 5.1 Ablation Study

We first conduct ablation studies to investigate the impact of different configurations. Please refer to



(a) MuST-C En⇒De



(b) MuST-C En⇒Es

Figure 3: StreamST performance of different methods.

## Appendix D.

### 5.2 Speech CoT Argumentation

StreamUni unifies streaming translation generation and policy decisions through speech CoT. The accuracy of outputs at each CoT stage significantly impacts overall StreamST performance, particularly under low-latency settings. To investigate this, we construct a low-latency speech evaluation dataset based on CoVoST2 En⇒Zh to assess generation capabilities across different CoT stages.

For dataset construction, we randomly truncate speech clip and obtain transcriptions using WhisperX (Bain et al., 2023), then generate reference translations using the DeepSeek-V3-0324 model (DeepSeek-AI et al., 2025b). We evaluate models trained with different schemes through speech CoT inference. More details are in Appendix E. As shown in Table 1, our approach achieves superior performance across all CoT stages, delivering excellent capabilities of policy decision and streaming translation generation.

### 5.3 Extending to Other LSLMs

Beyond the analytical experiments of our method, we further extend our evaluation to Qwen2.5-Omni-7B (Xu et al., 2025) to validate the generalizability

Training Settings	WER (↓)	SacreBLEU (↑)
Streaming CoT + Non-Streaming CoT	<b>20.74</b>	<b>35.22</b>
Non-Streaming CoT	27.83	33.60
Vanilla	31.62	33.34

Table 1: Performance of multiple stages of speech CoT under different training configurations. ‘Streaming CoT + Non-Streaming CoT’ denotes our employed training recipe. ‘Non-Streaming CoT’ only utilizes Non-Streaming CoT training data. ‘Vanilla’ represents the baseline without any further training.

Task	Base Model	LAAL(↓)	SacreBLEU(↑)
ST	Phi-4-Multimodal	N/A	28.55
	Qwen2.5-Omni	N/A	24.21
SimulST	Phi-4-Multimodal	1112.48	22.51
		1448.43	24.27
	Qwen2.5-Omni	949.36	20.64
		1449.83	21.80

Table 2: Performance of various vanilla LSLMs on ST and SimulST tasks. ‘ST’ denotes speech translation that utilizes complete speech inputs for translation, while ‘SimulST’ represents the simultaneous speech translation task that incorporates our proposed generation policy. The evaluation dataset is the MuST-C En⇒De sentence-level dataset.

of our approach across different LSLMs. The experimental results are presented in Table 2. Phi-4-Multimodal consistently outperforms Qwen-Omni on both ST and SimulST tasks, demonstrating that LSLMs with stronger speech translation capabilities achieve superior SimulST performance. This finding further validates that our StreamUni method can effectively leverage and scale with the enhanced capabilities of LSLMs, thereby demonstrating the generalizability of our approach.

## 6 Related Work

Streaming speech translation (StreamST) aims to generate real-time translations for continuously arriving speech stream, requiring the simultaneous completion of generation policy, segmentation policy, and streaming translation generation. Early research focused on sentence-level speech segments and is called simultaneous speech translation (SimulST), predominantly employing encoder-decoder architectures (Vaswani et al., 2017). Initial SimulST methods (Ma et al., 2020c) determine generation policy based on the number of input chunks. Subsequently, researchers explore content-

adaptive generation policy by leveraging auxiliary ASR tasks (Zeng et al., 2021; Chen et al., 2021; Zhang et al., 2024b), integrate-and-fire (Dong et al., 2022), monotonic attention (Communication et al., 2023), transducer (Liu et al., 2021; Tang et al., 2023), and CTC (Graves et al., 2006; Ma et al., 2023) to make decisions based on speech content. At the same time, some methods (Weller et al., 2021; Papi et al., 2023b; Omachi et al., 2022) attempt to accomplish translation by continuously refreshing the output translations.

With the advancement of Large Speech-Language Models (LSLMs), researchers have begun exploring their application to SimulST tasks (Agostinelli et al., 2024; Guo et al., 2025; Fu et al., 2025). Hibiki (Labiausse et al., 2025) even achieves simultaneous speech-to-speech translation in an end-to-end manner. However, relying solely on LSLMs for SimulST still requires coordination with multiple auxiliary models to achieve complete StreamST, introducing cascaded errors and hindering end-to-end optimization (Li et al., 2021). Consequently, researchers have attempted to develop unified methods capable of handling all StreamST tasks within a single model framework. Early attempts utilize attention mechanisms for generation and segmentation decisions (Papi et al., 2024), while subsequent work constructs dedicated policy-specific datasets to enable autoregressive prediction for policy decisions (Cheng et al., 2024; Ouyang et al., 2025). Nevertheless, these approaches suffer from significant challenges in large-scale data construction and advanced model transferability, while facing difficulties in fully leveraging the pre-training capabilities of foundation models.

## 7 Conclusion

In this paper, we propose StreamUni, a framework that efficiently enables unified LSLM to accomplish all subtasks of StreamST in a cohesive manner. By unifying different subtasks formats into autoregressive generation, StreamUni can achieve streaming translation with only a small amount of streaming CoT training data. Experiments show that our method efficiently attains state-of-the-art performance on StreamST tasks across multiple language directions with the same volume of training data. Furthermore, analytical experiments verify the effectiveness of each module in StreamUni as well as its practical usability.

648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
  
660  
661  
662  
663  
664  
665  
  
666  
667  
668  
669  
670  
671  
  
672  
673  
674  
  
675  
676  
677  
678  
679  
  
680  
681  
682  
683  
684  
  
685  
686  
687  
688  
689  
690  
691  
692  
693  
  
694  
695  
696  
697  
698  
699  
700  
701

## Limitations

StreamUni is designed to efficiently unify diverse tasks within the streaming speech translation framework. While our method achieves state-of-the-art (SOTA) performance across four language directions using only 400 hours of training data, the scale of the current dataset remains a primary constraint. We anticipate that scaling up the volume of training corpora and expanding the diversity of language pairs will yield further performance gains. Such explorations are reserved for future investigation.

## References

Victor Agostinelli, Max Wild, Matthew Raffel, Kazi Ahmed Asif Fuad, and Lizhong Chen. 2024. Simul-llm: A framework for exploring high-quality simultaneous translation with large language models. *arXiv preprint arXiv:2312.04691*.

Ashkan Alinejad and Anoop Sarkar. 2020. [Effectively pretraining a speech translation decoder with machine translation data](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8014–8020, Online. Association for Computational Linguistics.

Max Bain, Jaesung Huh, Tengda Han, and Andrew Senior. 2023. [Whisperx: Time-accurate speech transcription of long-form audio](#). *INTERSPEECH 2023*.

Junkun Chen, Mingbo Ma, Renjie Zheng, and Liang Huang. 2021. [Direct simultaneous speech-to-text translation assisted by synchronized streaming ASR](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*.

Shanbo Cheng, Zhichao Huang, Tom Ko, Hang Li, Ningxin Peng, Lu Xu, and Qini Zhang. 2024. [Towards achieving human parity on end-to-end simultaneous speech translation via llm agent](#). *Preprint*, arXiv:2407.21646.

Seamless Communication, Loïc Barrault, Yu-An Chung, Mariano Coria Meglioli, David Dale, Ning Dong, Mark Duppenthaler, Paul-Ambroise Duquenne, Brian Ellis, Hady Elsahar, Justin Haaheim, John Hoffman, Min-Jae Hwang, Hirofumi Inaguma, Christopher Klaiber, Iliia Kulikov, Pengwei Li, Daniel Licht, Jean Maillard, and 46 others. 2023. [Seamless: Multilingual expressive and streaming speech translation](#). *Preprint*, arXiv:2312.05187.

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, and 181 others. 2025a. [Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning](#). *Preprint*, arXiv:2501.12948.

DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, and 181 others. 2025b. [Deepseek-v3 technical report](#). *Preprint*, arXiv:2412.19437.

Mattia A. Di Gangi, Roldano Cattoni, Luisa Bentivogli, Matteo Negri, and Marco Turchi. 2019. [MuST-C: a Multilingual Speech Translation Corpus](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics*.

Qian Dong, Yaoming Zhu, Mingxuan Wang, and Lei Li. 2022. [Learning when to translate for streaming speech](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*.

Biao Fu, Donglei Yu, Minpeng Liao, Chengxi Li, Yidong Chen, Kai Fan, and Xiaodong Shi. 2025. [Efficient and adaptive simultaneous speech translation with fully unidirectional architecture](#). *Preprint*, arXiv:2504.11809.

Mattia A. Di Gangi, Matteo Negri, and Marco Turchi. 2019. [Adapting transformer to end-to-end spoken language translation](#). In *Interspeech 2019*, pages 1133–1137.

Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber. 2006. [Connectionist temporal classification: Labelling unsegmented sequence data with recurrent neural networks](#). In *ACM International Conference Proceeding Series*, pages 369–376.

Shoutao Guo, Shaolei Zhang, Zhengrui Ma, and Yang Feng. 2025. [Large language models are read/write policy-makers for simultaneous generation](#). *Preprint*, arXiv:2501.00868.

Zhichao Huang, Rong Ye, Tom Ko, Qianqian Dong, Shanbo Cheng, Mingxuan Wang, and Hang Li. 2023. [Speech translation with large language models: An industrial practice](#). *Preprint*, arXiv:2312.13585.

Javier Iranzo-Sánchez, Jorge Iranzo-Sánchez, Adrià Giménez, Jorge Civera, and Alfons Juan. 2024. [Segmentation-free streaming machine translation](#). *Preprint*, arXiv:2309.14823.

Tom Labiausse, Laurent Mazaré, Edouard Grave, Patrick Pérez, Alexandre Défossez, and Neil Zeghidour. 2025. [High-fidelity simultaneous speech-to-speech translation](#). *Preprint*, arXiv:2502.03382.

Ann Lee, Peng-Jen Chen, Changhan Wang, Jiatao Gu, Sravya Popuri, Xutai Ma, Adam Polyak, Yossi Adi, Qing He, Yun Tang, Juan Pino, and Wei-Ning Hsu. 2022. [Direct speech-to-speech translation with discrete units](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3327–3339, Dublin, Ireland. Association for Computational Linguistics.



869	Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2022. Robust speech recognition via large-scale weak supervision. <i>arXiv preprint arXiv:2212.04356</i> .	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> .	923
870			924
871			925
872			926
873	Ricardo Rei, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. 2022. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In <i>Proceedings of the Seventh Conference on Machine Translation (WMT)</i> , pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.	Dong Zhang, Shimin Li, Xin Zhang, Jun Zhan, Pengyu Wang, Yaqian Zhou, and Xipeng Qiu. 2023. SpeechGPT: Empowering large language models with intrinsic cross-modal conversational abilities. In <i>The 2023 Conference on Empirical Methods in Natural Language Processing</i> .	928
874			929
875			930
876			931
877			932
878			933
879			934
880			935
881	Yun Tang, Anna Sun, Hirofumi Inaguma, Xinyue Chen, Ning Dong, Xutai Ma, Paden Tomasello, and Juan Pino. 2023. Hybrid transducer and attention based encoder-decoder modeling for speech-to-text tasks. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8964–8986, Bangkok, Thailand. Association for Computational Linguistics.	Shaolei Zhang, Qingkai Fang, Shoutao Guo, Zhengrui Ma, Min Zhang, and Yang Feng. 2024a. Stream-Speech: Simultaneous speech-to-speech translation with multi-task learning. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8964–8986, Bangkok, Thailand. Association for Computational Linguistics.	936
882			937
883			938
884			939
885			940
886			941
887			942
888	Silero Team. 2024. Silero vad: pre-trained enterprise-grade voice activity detector (vad), number detector and language classifier. <a href="https://github.com/snakers4/silero-vad">https://github.com/snakers4/silero-vad</a> .	Shaolei Zhang, Qingkai Fang, Shoutao Guo, Zhengrui Ma, Min Zhang, and Yang Feng. 2024b. Stream-Speech: Simultaneous speech-to-speech translation with multi-task learning. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics</i> .	943
889			944
890			945
891			946
892	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In <i>Advances in Neural Information Processing Systems</i> , volume 30.	Shaolei Zhang and Yang Feng. 2023. End-to-end simultaneous speech translation with differentiable segmentation. In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> .	947
893			948
894			949
895			950
896			951
897	Changhan Wang, Anne Wu, and Juan Pino. 2020. Covost 2 and massively multilingual speech-to-text translation. <i>arXiv preprint 2007.10310</i> .	<b>A Truncation Policy in StreamUni</b>	952
898			953
899			954
900	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.	Our truncation policy is designed to truncate historical speech inputs and translations in real-time, enabling the model to focus on recent speech inputs while avoiding additional inference costs. To this end, we establish two key principles: (1) preventing elimination of speech segments containing untranslated content, which would compromise generation quality, and (2) avoiding removal of already-translated content from remaining speech segments, which would result in repetitive translations. Our approach first identifies the truncation timing for source speech inputs, then uses this as an anchor to determine the corresponding truncation point for output translations.	955
901			956
902			957
903			958
904			959
905	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models. In <i>Advances in Neural Information Processing Systems</i> .	For speech inputs, we consider appropriate truncation timing to be when users pause speaking or complete a sentence. Therefore, we design two triggering rules for speech truncation. The first rule targets prolonged user silence, while the second targets moments when users finish speaking a complete sentence. When neither condition is met for an extended period, causing the processed speech	960
906			961
907			962
908			963
909			964
910	Orion Weller, Matthias Sperber, Christian Gollan, and Joris Kluijvers. 2021. Streaming models for joint speech recognition and translation. <i>Preprint</i> , arXiv:2101.09149.		965
911			966
912			967
913			968
914	Jin Xu, Zhifang Guo, Jinzheng He, Hangrui Hu, Ting He, Shuai Bai, Keqin Chen, Jialin Wang, Yang Fan, Kai Dang, Bin Zhang, Xiong Wang, Yunfei Chu, and Junyang Lin. 2025. Qwen2.5-omni technical report. <i>Preprint</i> , arXiv:2503.20215.		969
915			970
916			971
917			972
918			973
919	Xi Xu, Wenda Xu, Siqi Ouyang, and Lei Li. 2024. Ca*: Addressing evaluation pitfalls in computation-aware latency for simultaneous speech translation. <i>Preprint</i> , arXiv:2410.16011.		974
920			
921			
922			

duration to exceed a predefined threshold (30 seconds), our method designates this moment as the truncation point.

After determining the speech truncation timing, we identify the corresponding truncation point for target translations to ensure semantic consistency between discarded content on both source and target sides. To achieve this alignment, we instruct the model to output all translations preceding the source truncation point and subsequently discard them.

This approach implements an effective truncation policy that maintains translation quality while ensuring computational efficiency.

## B Training and Evaluation Details

We provide comprehensive details of our training methodology. For training data construction, we focus on building streaming CoT data for the En $\Rightarrow$ Zh direction and incorporate an equal duration of non-streaming CoT data. For other language pairs, we directly utilize non-streaming data. The dataset released in this work is intended for academic research purposes only. Any commercial use is strictly prohibited. Our training details are detailed in Table 3.

For SimulST evaluation, we employ SimulEval (Ma et al., 2020a) as the standard assessment framework. For StreamST evaluation, we first utilize mWERSegmenter (Matusov et al., 2005) alignment tools to map the generated document-level translations to sentence-level references. Subsequently, we compute latency metrics and translation quality on the aligned sentences. We refer to these metrics as Stream LAAL (Papi et al., 2024) and Stream SacreBLEU, respectively.

## C Computation-Aware Latency

To validate the feasibility of our method for real-world deployment, we explicitly incorporate computational latency into our evaluation framework when conducting experiments on the NVIDIA GeForce RTX 3090 (a consumer-grade GPU). Specifically, we adopt the computation-aware LAAL metric (Papi et al., 2022b; Xu et al., 2024), which quantifies the end-to-end latency from the user’s speech input to the model’s translation output. The table below reports our method’s StreamST performance on the MuST-C En $\Rightarrow$ De and En $\Rightarrow$ Es datasets, with both computation-aware

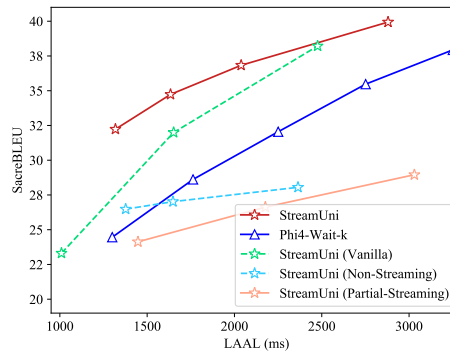


Figure 4: The SimulST performance on En $\Rightarrow$ Zh with different policies and training methods: Vanilla (no extra training), Non-Streaming (trained solely on non-streaming data), and Partial-Streaming (encourages the model to predict the corresponding translation based on the partial speech input, rather than translation corresponding to the unsegmented speech input).

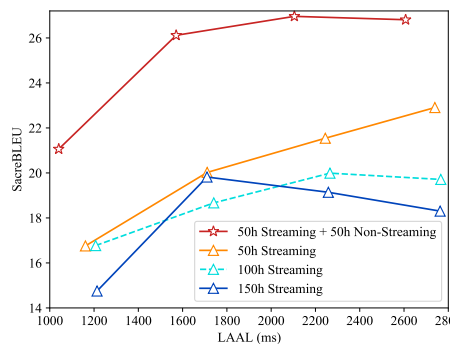


Figure 5: The SimulST performance of different training data recipes on MuST-C En $\Rightarrow$ De dataset. Our method utilizes 50 hours each of proposed streaming and non-streaming data, whereas other methods only employ the proposed streaming data with varying total durations.

latency and non-computation-sensitive metrics included for comprehensive comparison.

As can be seen from the Table 4 and Table 5, our method achieves promising performance with a latency of approximately 3 seconds when computational costs are taken into account, and around 2 seconds when computational delays are not considered. Notably, these results are obtained without leveraging any inference optimization frameworks or advanced GPUs. It is anticipated that the adoption of the aforementioned optimization techniques will enable us to achieve even much lower latency.

## D Ablation Study

Figure 4 presents a performance comparison of our method under various training methods and gen-

Hyperparameters			Settings
	Base_model	Base_model	Phi-4-Multimodal
LSLM	Training Details	batch_size	32
		learning_rate	4e-5
		weight_decay	0.01
		lr_scheduler	WarmupLR
		betas	(0.9, 0.95)
		optimizer	AdamW
		zero_optimization	stage_2

Table 3: Settings of StreamUni.

Computation-Aware Stream LAAL (ms)	Stream LAAL (ms)	Stream SacreBLEU
3193.94	1762.85	25.65
3690.02	2223.69	26.89

Table 4: Computation-Aware latency and translation results on MuST-C En $\Rightarrow$ De task.

eration policies. Unlike Phi4-Wait- $k$ , which employs heuristic rules for generation decisions without considering speech content (Ma et al., 2019), our method determines generation timing by detecting valid speech inputs, thereby achieving superior performance through more informed generation policies. Beyond generation policy, our proposed streaming CoT training scheme enhances performance across all latency settings, particularly under low-latency. However, the streaming CoT training data must be combined with non-streaming data to achieve maximum performance gains.

To validate this hypothesis, we conduct experiments using training data from a single language direction. As illustrated in Figure 5, simply increasing data volume when using only streaming CoT data fails to yield performance improvements. Superior performance across different latency settings is achieved exclusively when both streaming and non-streaming CoT data are employed simultaneously. We hypothesize that this mixed-data approach effectively stimulates the streaming generation capabilities of model while enabling it to perceive complete speech input boundaries, thereby preventing over-translation and achieving enhanced overall performance.

Furthermore, we investigate the effectiveness of our proposed truncation policy. Rather than comparing the accuracy of model-determined truncation timing against official truncation points, we focus on the final generation quality, which represents our ultimate objective. Given document-level

speech inputs with an average duration exceeding 10 minutes (Di Gangi et al., 2019), we evaluate the generation quality of our fine-tuned model under different truncation policies. Table 6 illustrates the results, where we report document-level metrics rather than sentence-level metrics after alignment. Notably, while our proposed truncation strategy performs slightly below the human-annotated policy on SacreBLEU, it surpasses official annotation on COMET score. This demonstrates the effectiveness of our approach and provides valuable insights for future research utilizing semantic alignment models to implement truncation policies.

Beyond that, We also provide an analysis of why our method, despite a slightly lower SacreBLEU, outperforms the officially provided truncation policy in terms of COMET. We conduct a detailed analysis of the average duration of processed speech inputs and average translation length obtained from two truncation policies. The results are shown in the Table 7. We find that the truncation policy of our method enables the model to process longer speech segments and reference lengthier historically generated translations. This provides sufficient context for the model during translation, thus ensuring that our truncation policy outperforms the truncation approach of MuST-C.

Furthermore, we also demonstrate the impact of adopting Speech CoT on translation performance. Herein, we use the RealSI dataset (Cheng et al., 2024) to explore full-sentence speech translation performance with and without Speech CoT. Re-

Computation-Aware Stream LAAL (ms)	Stream LAAL (ms)	Stream SacreBLEU
3077.07	1802.59	29.02
3405.03	2087.39	30.18

Table 5: Computation-Aware latency and translation results on MuST-C En⇒Es task.

Direction	Truncation	COMET	SacreBLEU
En⇒De	official	82.45	<b>32.51</b>
	Model	<b>83.42</b>	31.59
En⇒Es	official	80.83	<b>35.84</b>
	Model	<b>82.86</b>	34.97

Table 6: Performance on MuST-C with document-level speech inputs when the generation policy is disabled and only truncation is employed. ‘Human’ uses annotated truncation timing, while ‘Model’ lets the system decide automatically.

alSI (Cheng et al., 2024) is a speech test set from real-world scenarios containing English⇒Chinese directions. It can be observed in Table 8 that CoT also brings significant improvements in translation performance.

## E Speech CoT Argumentation

For streaming speech translation, the key challenge lies in real-time performance. This challenge is amplified under extremely low latency, where very short input segments make accurate translation and policy decision-making especially important. To evaluate these aspects, we construct a dedicated Low-Latency Speech Evaluation dataset. This evaluation set was derived from conventional speech translation corpora (consisting of speech segments, transcripts, and translations) through the following modifications:

- For each complete speech segment, we randomly sample speech prefixes with shorter duration based on the given complete speech segment.
- Using WhisperX (Bain et al., 2023), we obtain word-level timestamps for the complete speech segment, which allows us to extract the corresponding ground-truth transcript prefixes for the sampled speech prefixes.
- With DeepSeek-V3-0324 (DeepSeek-AI et al., 2025b), we generate ground-truth translations of the transcript prefixes, yielding the low-latency speech evaluation dataset.

During evaluation, the model is encouraged to generate intermediate results of speech CoT based on speech prefixes at different stages. For the ASR transcription outputs of speech CoT, we compute the WER against the ground-truth transcription prefixes to assess its low-latency transcription capability, which further reflects its policy-decision ability under low latency for our StreamUni. For the translation results of speech CoT, we calculate SacreBLEU against the ground-truth references to measure its low-latency translation capability. Based on the constructed evaluation dataset, We evaluate models trained with different schemes through speech CoT inference.

The detailed experimental results are shown in Table 1. Our employed ‘Streaming CoT + Non-Streaming CoT’ training scheme achieves lower WER and higher SacreBLEU scores, our approach achieves superior performance across all CoT stages, delivering excellent capabilities of policy decision and streaming translation generation.

<b>Settings</b>	<b>Avg Duration</b>	<b>Avg Translation Length</b>
StreamUni	9.34	27.91
MuST-C	5.77	15.6

Table 7: Analysis of input and output patterns of different truncation policies on MuST-C En $\Rightarrow$ De task.

<b>Settings</b>	<b>SacreBLEU</b>
Direct Trans	22.08
CoT Trans	25.23

Table 8: Full-Sentence Translation Performance: With vs. Without Speech CoT Argumentation.