

# Adapting Offline Speech Translation Models for Streaming with Future-Aware Distillation and Inference

Anonymous ACL submission

## Abstract

A popular approach to streaming speech translation is to employ a single offline model with a *wait-k* policy to support different latency requirements, which is simpler than training multiple online models with different latency constraints. However, there is a mismatch problem in using a model trained with complete utterances for streaming inference with partial input. We demonstrate that speech representations extracted at the end of a streaming input are significantly different from those extracted from a complete utterance. To address this issue, we propose a new approach called Future-Aware Streaming Translation (FAST) that adapts an offline ST model for streaming input. FAST includes a Future-Aware Inference (FAI) strategy that incorporates future context through a trainable masked embedding, and a Future-Aware Distillation (FAD) framework that transfers future context from an approximation of full speech to streaming input. Our experiments on the MuST-C EnDe, EnEs, and EnFr benchmarks show that FAST achieves better trade-offs between translation quality and latency than strong baselines. Extensive analyses suggest that our methods effectively alleviate the aforementioned mismatch problem between offline training and online inference.

## 1 Introduction

Streaming speech translation (ST) systems generate real-time translations by incrementally processing audio frames, unlike their offline counterparts that have access to complete utterances before translating. Typically, streaming ST models use unidirectional encoders (Ren et al., 2020; Ma et al., 2020b; Zeng et al., 2021) and are trained with a read/write policy that determines whether to wait for more speech frames or emit target tokens. However, it can be expensive to maintain multiple models to satisfy different latency requirements (Zhang and Feng, 2021; Liu et al., 2021a) in real-world applications. Recently, some works (Papi et al.,

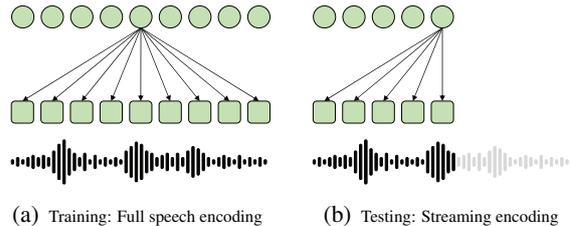


Figure 1: (a) and (b) represent the input mismatch between offline training and streaming testing.

2022; Dong et al., 2022) have shown that a single offline models with bidirectional encoders (such as Wav2Vec2.0 (Baevski et al., 2020)) can be adapted to streaming scenarios with a *wait-k* policy (Ma et al., 2019) to meet different latency requirements and achieve comparable or better performance, partially due to the more powerful bidirectional encoders. However, there is an inherent mismatch in using a model trained with complete utterances on incomplete streaming speech during online inference (Ma et al., 2019).

Intuitively, speech representations extracted from streaming inputs (Figure 1(b)) are less informative than those from full speech encoding (Figure 1(a)) due to limited future context, especially toward the end of the streaming inputs, which can be exacerbated by the aforementioned mismatch problem. This raises a natural question: how much do the speech representations differ between the two inference modes? We analyze the gap in speech representations, measured by cosine similarity, at different positions in the streaming input compared to using the full speech (Section 3). We observe a significantly greater gap for representations closer to the end of a streaming segment, with an average similarity score as low as 0.2 for the last frame, and the gap quickly narrows for earlier frames (Figure 2). Additionally, we observe more degradation in translation quality for utterances with the greatest gap in speech representations between online

073 and offline inference (see Appendix B.2).

074 We conjecture that the lack of future contexts  
075 at the end of streaming inputs can be detrimental  
076 to streaming speech translation when using an  
077 offline model. To this end, we propose a novel  
078 Future-Aware Inference (FAI) strategy. This ap-  
079 proach is inspired by masked language models’  
080 ability (Baevski et al., 2020) to construct repre-  
081 sentations for masked tokens from their context.  
082 Specifically, we append a few mask embeddings  
083 to the end of the streaming input and leverage the  
084 acoustic encoder (Wav2Vec2.0)’s ability to implic-  
085 itly construct representations for future contexts,  
086 which can lead to more accurate representations for  
087 the other frames in the streaming input.

088 Furthermore, we propose a Future-Aware Dis-  
089 tillation (FAD) framework that adapts the offline  
090 model to extract representations from streaming  
091 inputs that more closely resemble those from full  
092 speech encoding. We expand the original stream-  
093 ing input with two types of future contexts: one  
094 with  $m$  oracle speech tokens for the teacher model,  
095 and another with  $m$  mask tokens for the student  
096 model, which is initialized from the teacher model.  
097 We minimize several distillation losses between  
098 the output of the teacher and student models. By  
099 incorporating additional oracle future contexts, the  
100 speech representations for the frames in the origi-  
101 nal streaming input extracted by the teacher model  
102 resemble those when the full speech is available.  
103 FAD aims to adjust the offline model to extract  
104 similar representations for streaming input as it  
105 would for full speech. In combination with FAI,  
106 we improve the model’s ability to extract quality  
107 representations during both training and inference,  
108 alleviating the aforementioned mismatch problem.  
109 We refer to our approach as FAST, which stands  
110 for Future-Aware Streaming Translation.

111 We conducted experiments on the MuST-C  
112 EnDe, EnEs, and EnFr benchmarks. The results  
113 show that our methods outperform several strong  
114 baselines in terms of the trade-off between transla-  
115 tion quality and latency. Particularly, in the lower  
116 latency range (when AL is less than 1000ms), our  
117 approach achieved BLEU improvements of 12 in  
118 EnDE, 16 in EnEs, and 14 in EnFr over baseline.  
119 Extensive analyses demonstrate that our future-  
120 aware approach significantly reduces the represen-  
121 tation gap between partial streaming encoding and  
122 full speech encoding.

## 2 Background and Related Work

Speech translation systems can be roughly catego-  
rized into non-streaming (offline) and streaming  
(online) depending on the inference mode. Re-  
gardless of the inference mode, speech transla-  
tion models typically employ the encoder-decoder  
architecture and are trained on an ST corpus  
 $\mathcal{D} = \{(\mathbf{x}, \mathbf{z}, \mathbf{y})\}$ , where  $\mathbf{x} = (x_1, \dots, x_T)$  de-  
notes an audio sequence,  $\mathbf{z} = (z_1, \dots, z_I)$  and  
 $\mathbf{y} = (y_1, \dots, y_J)$  the corresponding source tran-  
scription and target translation respectively.

**Non-Streaming Speech Translation** For the  
non-streaming ST task, the encoder maps the en-  
tire input audio  $\mathbf{x}$  to the speech representations  
 $\mathbf{h}$ , and the decoder generates the  $j$ -th target to-  
ken  $y_j$  conditional on the full representations  $\mathbf{h}$   
and the previously generated tokens  $y_{<j}$ . The de-  
coding process of non-streaming ST is defined as  
 $p(\mathbf{y} | \mathbf{x}) = \prod_{j=1}^J p(y_j | \mathbf{x}, \mathbf{y}_{<j})$ .

A significant amount of works have focused on  
non-streaming ST, including pre-training (Wang  
et al., 2020; Dong et al., 2021a; Tang et al., 2022;  
Ao et al., 2022), multi-task learning (Liu et al.,  
2020; Indurthi et al., 2020, 2021), data augmenta-  
tion (Pino et al., 2019; Di Gangi et al., 2019b; Mc-  
Carthy et al., 2020), knowledge distillation (Dong  
et al., 2021b; Zhao et al., 2021; Du et al., 2022),  
and cross-modality representation learning (Tang  
et al., 2021; Fang et al., 2022; Ye et al., 2022).

**Streaming Speech Translation** A streaming ST  
model generates the  $j$ -th target token  $y_j$  based  
on streaming audio prefix  $\mathbf{x}_{\leq g(j)}$  and the pre-  
vious tokens  $y_{<j}$ , where  $g(j)$  is a monotonic  
non-decreasing function representing the ending  
timestamp of the audio prefix that needs to be  
consumed to generate the  $j$ -th word. The de-  
coding probability is calculated as  $p(\mathbf{y} | \mathbf{x}) =$   
 $\prod_{j=1}^J p(y_j | \mathbf{x}_{\leq g(j)}, \mathbf{y}_{<j})$ .

Thus, a streaming ST model requires a policy to  
determine whether to wait for more source speech  
or emit new target tokens. Recent studies (Ma et al.,  
2020b; Ren et al., 2020; Zeng et al., 2021; Dong  
et al., 2022) make read/write decisions based on a  
variant of the *wait- $k$*  policy (Ma et al., 2019) that  
was initially proposed for streaming text transla-  
tion, which alternates write and read operations  
after reading the first  $k$  source tokens. Because  
there is no explicit word boundaries in a stream-  
ing audio, several works attempt to detect word  
boundaries in the audio sequence by fixed length  
(Ma et al., 2020b), Connectionist Temporal Classi-

174 fication (Ren et al., 2020; Zeng et al., 2021; Papi  
 175 et al., 2022), ASR outputs (Chen et al., 2021), or  
 176 continuous-integrate-and fire (Dong et al., 2022;  
 177 Chang and yi Lee, 2022). Moreover, some studies  
 178 (Arivazhagan et al., 2019; Ma et al., 2020c; Zhang  
 179 et al., 2020; Schneider and Waibel, 2020; Miao  
 180 et al., 2021; Zhang and Feng, 2022a,c; Zhang et al.,  
 181 2022; Chang and yi Lee, 2022; Liu et al., 2021b;  
 182 Zhang and Feng, 2022b) explore adaptive policies  
 183 to dynamically decide when to read or write for  
 184 streaming text and/or streaming speech translation.  
 185 Zhang and Feng (2022d) fill future source posi-  
 186 tions with positional encoding to introduce future  
 187 information during training for simultaneous ma-  
 188 chine translation (MT) within the prefix-to-prefix  
 189 framework. In this paper, we focus on a matter  
 190 less attended to – how to alleviate the mismatch  
 191 between offline training and online inference.

192 **Knowledge Distillation for Streaming Trans-**  
 193 **lation** Existing studies on streaming text and/or  
 194 speech translation usually introduce future informa-  
 195 tion by distilling sequence-level knowledge from  
 196 offline MT (Ren et al., 2020; Zhang et al., 2021;  
 197 Liu et al., 2021b; Zhu et al., 2022; Deng et al.,  
 198 2023) and online MT (Zaidi et al., 2021). More-  
 199 over, Ren et al. (2020) leverage the knowledge from  
 200 the multiplication of attention weights matrices of  
 201 streaming ASR and MT models to supervise the  
 202 attention of the streaming ST model. However, our  
 203 FAD aims to reduce the representation gap between  
 204 full speech and streaming speech.

### 205 3 Preliminary Analysis

206 In this section, we examine the mismatch problem  
 207 in Transformer-based (Vaswani et al., 2017) ST  
 208 architecture between offline training and online de-  
 209 coding. In offline full-sentence ST, the speech rep-  
 210 resentation of each frame is obtained by attending  
 211 to all frames, including future frames, in the trans-  
 212 former encoder layers. Recently, a common ap-  
 213 proach in speech translation is to stack a pre-trained  
 214 Wav2Vec2.0 (Baevski et al., 2020) as the acoustic  
 215 encoder with a semantic MT encoder-decoder, re-  
 216 sulting in state-of-the-art performance in the ST  
 217 task (Han et al., 2021; Dong et al., 2022; Fang  
 218 et al., 2022; Ye et al., 2022). This approach lever-  
 219 ages the ability of Wav2Vec2.0 pre-training to learn  
 220 better speech representations.

221 When applying an offline model to streaming  
 222 inference, the lack of future frames causes an ap-  
 223 parent mismatch problem, which can lead to a de-

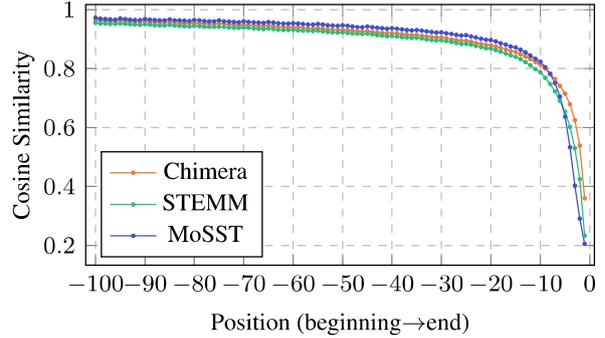


Figure 2: The average cosine similarity  $\bar{s}_{-\tau}$  of the end 100 positions in the streaming speech.

224 terioration in the extracted speech representations.  
 225 To quantify this effect, we examine three offline  
 226 ST models trained on the MuST-C EnDe dataset  
 227 using the Chimera (Han et al., 2021), STEM  
 228 (Fang et al., 2022), and MoSST (Dong et al., 2022)  
 229 architectures, with a trainable acoustic encoder ini-  
 230 tialized from Wav2Vec2.0. We conduct analysis  
 231 on the tst-COMMON set with a duration between  
 232 2s and 10s by removing outliers and noisy data,  
 233 resulting 1829 examples.

234 For an input sequence of audio frames  $\mathbf{x} =$   
 235  $(x_1, \dots, x_T)$ , the convolutional subsampler of  
 236 Wav2Vec2.0 shrinks the length of the raw audio  
 237 by a factor 320 and outputs the full speech rep-  
 238 resentation sequence  $\mathbf{a}$ . For readability reasons,  
 239 we uniformly use the notation  $T$  to denote the se-  
 240 quence length of  $\mathbf{a} = (a_1, \dots, a_T)$ . This simplified  
 241 notation does not undermine any of our conclusions  
 242 while making the equations for readable. For stream-  
 243 ing input  $\forall t \leq T$ ,  $\hat{\mathbf{x}}_t = (x_1, \dots, x_t)$ , Wav2Vec2.0  
 244 will output the representation  $\hat{\mathbf{a}}_t = (\hat{a}_{t,1}, \dots, \hat{a}_{t,t})$ .

245 To quantify the difference in speech representa-  
 246 tions between offline and online inputs, we com-  
 247 pute the cosine similarity  $s_{t,t'}$  between the speech  
 248 representation at the  $t'$ -th ( $t' \leq t$ ) position in the  
 249 streaming audio input  $\hat{\mathbf{x}}_t$  and at the same position  
 250 with full-sentence encoding. We then calculate the  
 251 statistics  $\bar{s}_{-\tau}$  by averaging the cosine similarity  
 252 over both the testset  $\mathcal{B}$  and the time dimension with  
 253 a reverse index  $-\tau$  corresponding to a position  
 254  $\tau - 1$  frames before the end of the streaming input.

$$255 s_{t,t'}(\mathbf{x}) = \cos(\hat{a}_{t,t'}, a_{t'}), \forall t' \leq t, \quad (1)$$

$$256 \bar{s}_{-\tau} = \frac{1}{|\mathcal{B}|} \sum_{\mathbf{x} \in \mathcal{B}} \frac{1}{|\mathbf{x}| - \tau + 1} \sum_{t=\tau}^{|\mathbf{x}|} s_{t,t-\tau+1}(\mathbf{x}) \quad (2)$$

257 Figure 2 displays the  $\bar{s}_{-\tau}$  curve for the last 100  
 258 positions in streaming inputs. For  $\tau > 10$ , the

259 averaged cosine similarity  $\bar{s}_{-\tau}$  is greater than 0.8, 308  
 260 indicating that the representations at those posi-  
 261 tions in a streaming input are similar to those with  
 262 the full speech. However, the curve shows a sharp  
 263 decline in the averaged cosine similarity  $\bar{s}_{-\tau}$   
 264 for the ending positions, particularly for the last one  
 265 ( $\tau = 1$ ), suggesting that the mismatch problem can  
 266 significantly affect the quality of speech represen-  
 267 tation for these positions. We provide additional  
 268 analysis in Appendix B.

## 269 4 Method

270 To address the mismatch problem between of-  
 271 fline training and online inference, we propose a  
 272 novel methodology called Future-Aware Stream-  
 273 ing Translation (FAST). This approach adapts an  
 274 offline ST model for streaming scenarios by us-  
 275 ing a Future-Aware Inference (FAI) strategy during  
 276 inference and a Future-Aware Distillation (FAD)  
 277 strategy during training. An overview of our pro-  
 278 posed method is depicted in Figure 3.

### 279 4.1 Model Architecture

280 Unlike previous works (Ren et al., 2020; Ma et al.,  
 281 2020b; Zeng et al., 2021; Liu et al., 2021a) that  
 282 require training multiple streaming models for dif-  
 283 ferent latency requirements, our goal is to train one  
 284 single offline model to meet the requirements. The  
 285 overall architecture depicted in Figure 3(a) consists  
 286 of an acoustic encoder, an acoustic boundary detec-  
 287 tor, a semantic encoder, and a translation decoder.  
 288 **Acoustic encoder:** The pre-trained Wav2Vec2.0  
 289 is adopted as the acoustic encoder to learn a better  
 290 speech representation (Ye et al., 2021, 2022).

291 **Acoustic boundary detector:** To enable the of-  
 292 fline ST model to perform chunk-wise streaming  
 293 inference, we use a Continuous Integrate-and-Fire  
 294 (CIF) module (Dong and Xu, 2020) as the acoustic  
 295 boundary detector to dynamically locate the acous-  
 296 tic boundaries of speech segments following (Yi  
 297 et al., 2021; Dong et al., 2022). The CIF module  
 298 generates an integration weight  $\alpha_t$  for each acous-  
 299 tic representation  $a_t$  by Wav2Vec2.0. Then, CIF  
 300 accumulates  $\alpha_t$  in a step-by-step way. When the  
 301 accumulation reaches a certain threshold (e.g. 1.0),  
 302 the acoustic representations corresponding to these  
 303 weights are integrated into a single hidden represen-  
 304 tation  $h_j$  by weighted average, indicating a found  
 305 token boundary. The shrunk representations  $\mathbf{h}$  will  
 306 be fed into the semantic encoder. To learn the cor-  
 307 rect acoustic boundaries, we use the source text

length  $J$  as the weakly supervised signal.

$$260 \mathcal{L}_{\text{CIF}} = \left\| J - \sum_{t=1}^T \alpha_t \right\|_2 \quad (3) \quad 309$$

310 There are two benefits of using CIF as a boundary  
 311 detector. For offline ST model, it can address the  
 312 length gap between speech and text. It can also pro-  
 313 vide the acoustic boundaries to perform read/write  
 314 policies for streaming inference.

315 **Semantic encoder and Translation decoder:** The  
 316 standard transformer (Vaswani et al., 2017) com-  
 317 posed of  $L_e$  encoder layers and  $L_d$  decoder layers  
 318 is used. The translation loss is defined as:

$$272 \mathcal{L}_{\text{ST}}(\mathbf{x}, \mathbf{y}) = - \sum_{j=1}^J \log p(y_j | y_{<j}, \mathbf{x}) \quad (4) \quad 319$$

### 320 4.2 Future-Aware Inference

321 The offline ST model is trained with the following  
 322 objective function:

$$278 \mathcal{L}_{\text{offline}} = \mathcal{L}_{\text{ST}} + \lambda \cdot \mathcal{L}_{\text{CIF}} \quad (5) \quad 323$$

324 where  $\lambda$  is a hyper-parameter to balance two losses.

325 Based on the analysis in Section 3, we find that  
 326 it is only necessary for the offline ST model to be  
 327 aware of a short future during streaming encoding.  
 328 Thus, we first propose a Future-Aware Inference  
 329 (FAI) strategy to enhance the representations of  
 330 streaming speech in Figure 3(b).

331 In this strategy, the streaming inference is di-  
 332 rectly performed on offline ST model without fine-  
 333 tuning. Particularly, we use the mask tokens of  
 334 Wave2Vec2.0 as the pseudo future context and ap-  
 335 pend them to the speech tokens generated from  
 336 the already consumed speech frames. Because  
 337 the mask token embedding is trainable when pre-  
 338 training Wave2Vec2.0, and the contrastive loss is to  
 339 identify the quantized latent audio representation of  
 340 masked regions based on unmasked context, this is  
 341 intuition that mask tokens can possibly encode fu-  
 342 ture context. In addition, the masking strategy dur-  
 343 ing pre-training results in approximately 49% of all  
 344 time steps being masked with a mean span length  
 345 of 300ms, it also guarantees that Wav2vec2.0 is  
 346 able to extract better speech representations even  
 347 with the presence of large amount of mask tokens.

348 Wav2Vec2.0 consists of a multi-layer convolu-  
 349 tional subsampler  $f_c$  and a Transformer encoder  $f_e$ .  
 350 During our online inference, for each audio prefix  
 351  $\hat{\mathbf{x}}_t = (x_1, \dots, x_t)$ , the  $f_c$  first outputs streaming  
 352 speech tokens  $\hat{\mathbf{c}}_t = (c_1, \dots, c_\tau)$ , where  $\hat{\mathbf{c}} \in \mathbb{R}^{\tau \times d}$   
 353 and  $d$  is the dimension of model and  $\tau$  is the

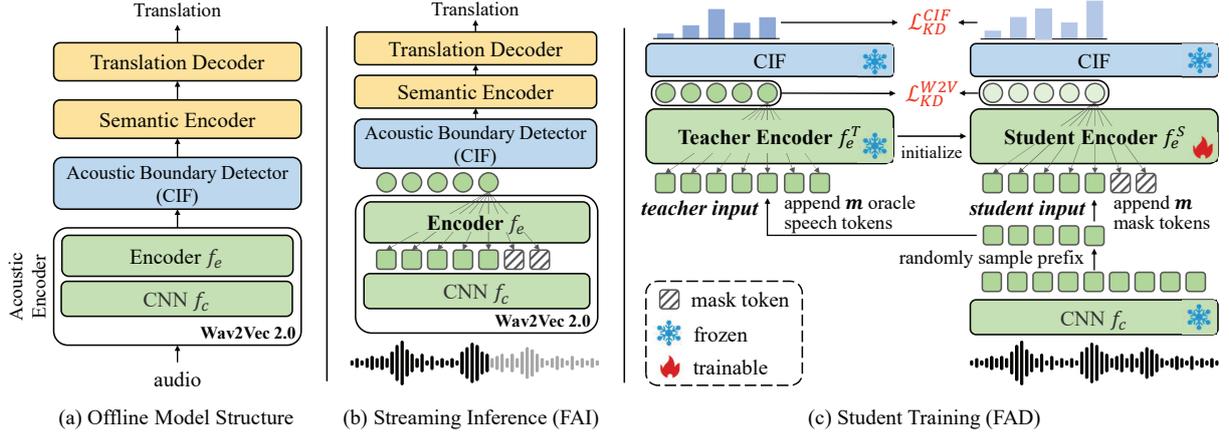


Figure 3: Illustration of offline ST model and proposed methods FAI and FAD.

sequence length after convolutional subsampling. Then, we concatenate the streaming speech tokens  $\hat{\mathbf{c}}$  and  $m$  mask token embeddings  $\mathbf{e} \in \mathbb{R}^d$  along the time dimension, resulting in a longer sequence of speech tokens  $\in \mathbb{R}^{(\tau+m) \times d}$ . The new speech tokens are then fed into the Transformer encoder  $f_e$ , but only the first  $\tau$  encoder outputs (i.e., speech features) will be kept for the CIF module because, as discussed in Section 3, the last  $m$  speech features are of poor quality and adversely affect translation quality. Then, if an acoustic boundary is detected by the CIF module, the decoder will emit new words based on *wait-k* policy, otherwise, the streaming speech continues to be read. The FAI strategy is outlined in Algorithm 1 in Appendix.

### 4.3 Future-Aware Distillation

Even FAI considers mask tokens as the pseudo future context, it is still preferred to leverage the future oracle speech tokens, which is unavailable during inference. Therefore, we take one step further by proposing a fine-tuning method – Future-Aware Distillation (FAD). It aims to distill the knowledge from teachers with oracle future contexts into students with pseudo future contexts.

The **teacher** model is the offline ST by optimizing Eq. (5) and is frozen. The **student** model has exactly the same architecture as the teacher and is initialised from the teacher. However, the semantic encoder and translation decoder are frozen to retrain offline-trained ST performance.

**Training** A naive solution is to distill knowledge from the full speech into every possible streaming speech for each audio. However, since the length of speech tokens is typically very large, e.g., 300 on average, it will be computational prohibitive. To

this end, we propose a simple and efficient implementation via random sampling.

Given a full audio waveform  $\mathbf{x}$ ,  $f_c$  outputs the speech tokens  $\mathbf{c} \in \mathbb{R}^{T \times d}$ . We randomly sample an integer  $t \in [1, T]$  to construct the streaming speech token  $\mathbf{c}_{\leq t}$ . Then, we define the teacher input of  $f_e$  with oracle future context as following:

$$\hat{\mathbf{c}}^T = \mathbf{c}_{1:t+m} \in \mathbb{R}^{(t+m) \times d}, \quad (6)$$

where  $m$  is a hyper-parameter to denote the number of future contexts. The most straightforward approach is to use the full speech as the teacher input. However, due to the bidirectional acoustic encoder, the streaming speech representation of the same position constantly changes when consuming new frames.

To maintain consistency with the inference method FAI, we use the mask tokens as the pseudo future context and append them to the sampled speech tokens to construct the student input.

$$\hat{\mathbf{c}}^S = \text{Concat}\{\mathbf{c}_{1:t}; m \times [\mathbf{e}]\} \in \mathbb{R}^{(t+m) \times d}, \quad (7)$$

where  $\mathbf{e} \in \mathbb{R}^d$  is the mask embedding.

We can obtain the streaming speech representations from teacher  $f_e^T$  and student  $f_e^S$ . Then the first  $t$  speech representations are fed into the CIF module to derive the teacher and student weight sequence. Concretely, they can be written as follows.

$$\hat{\mathbf{a}}^T, \hat{\mathbf{a}}^S = f_e^T(\hat{\mathbf{c}}^T), f_e^S(\hat{\mathbf{c}}^S) \quad (8)$$

$$\alpha_{1:t}^T, \alpha_{1:t}^S = \text{CIF}(\hat{\mathbf{a}}_{1:t}^T), \text{CIF}(\hat{\mathbf{a}}_{1:t}^S) \quad (9)$$

Eventually, two distillation losses are proposed to reduce the speech representation gap.

$$\mathcal{L}_{KD}^{W2V} = 1 - \text{cosine}(\hat{\mathbf{a}}_{1:t}^S, \hat{\mathbf{a}}_{1:t}^T) \quad (10)$$

$$\mathcal{L}_{KD}^{\text{CIF}} = \sum_{\tau=1}^t \text{KL}(\alpha_{\tau}^T \parallel \alpha_{\tau}^S) \quad (11)$$

The first loss is to directly minimize the streaming speech representations with cosine similarity. The second loss is to learn more correct acoustic boundaries for online inference by calculating the KL-divergence between two weight distributions. Note that according previous analysis in Sec 3, the representations of the first  $t$  speech tokens after  $f_e^T$  should have high quality if  $m > 10$ , so only the first  $t$  speech representations are taken into account for loss calculation.

**Optimization** The total training objective of the FAD can be written as  $\mathcal{L} = \mathcal{L}_{KD}^{W2V} + \mathcal{L}_{KD}^{CIF}$ . The overall training procedure of the proposed method is shown in Figure 3(c).

## 5 Experiments

### 5.1 Experimental Settings

**Datasets** We evaluate our approach on MuST-C V1 English-German (EnDe), English-Spanish (EnEs) and English-French (EnFr) datasets (Di Gangi et al., 2019a), where limited previous works discussed the En-Fr streaming ST with BLEU-latency curve. All the corpora contain source audios, source transcriptions, and target translations, and the results reported are conducted on the corresponding tst-COMMON set. For speech data, we normalize the raw audio wave to the range of  $[-1, 1)$ . For text data, we keep punctuation and remove non-printing characters, and remain case-sensitive. For each translation direction, the unigram Sentence-Piece<sup>1</sup> model (Kudo and Richardson, 2018) is used to learn a shared vocabulary of size 10k.

**Model Configuration** For the acoustic encoder, we use Wav2vec2.0<sup>2</sup> (Baevski et al., 2020) following the base configurations. We construct the acoustic boundary detector by applying the CIF (Yi et al., 2021) on the last dimension of speech representation. We use 8 and 6 layers for the semantic encoder and the translation decoder respectively, with 4 attention heads and 768 hidden units.

**Training** The detailed training schedule of the offline ST model can refer to Appendix C. We set the length  $m$  of future context tokens to 50 for both FAD and FAI. All hyper-parameters are tuned on EnDe devset and applied to other language pairs. We train all models with 3.2 million frames per batch on 8 Nvidia Tesla V100 GPUs. We imple-

<sup>1</sup><https://github.com/google/sentencepiece>

<sup>2</sup>[https://dl.fbaipublicfiles.com/fairseq/wav2vec/wav2vec\\_small.pt](https://dl.fbaipublicfiles.com/fairseq/wav2vec/wav2vec_small.pt)

ment our models with Fairseq<sup>3</sup> (Ott et al., 2019). **Inference** We average the checkpoints of the best 10 epochs on development set for evaluation. We perform streaming-testing with the *wait-k* policy.  $k$  is counted by the detected acoustic units from the CIF module. To follow the tradition in simultaneous translation (Zeng et al., 2021; Dong et al., 2022), we do not rewrite the tokens that have already been generated.

**Evaluation Metrics** We use SacreBLEU<sup>4</sup> for the translation quality. The latency is evaluated with Average Latency (AL) (Ma et al., 2019), Average Proportion (AP) (Cho and Esipova, 2016), and Differentiable Average Lagging (DAL) (Cherry and Foster, 2019) in the SimulEval<sup>5</sup> (Ma et al., 2020a).

**System Settings** We compare our method with several strong end-to-end streaming ST approaches. (i) *SimulSpeech* (Ren et al., 2020) and *RealTranS* (Zeng et al., 2021) use uni-directional encoder rather than bidirectional one. (ii) *MoSST* (Dong et al., 2022) applies an offline-trained model with a monotonic segmentation module for streaming testing and achieves competitive performance. (iii) *MMA-SLM* (Indurthi et al., 2022) enhances monotonic attention to make better read/write decisions by integrating future information from language models. (iv) *ITST* (Zhang and Feng, 2022b) learns an adaptive read/write policy by quantifying the transported information weight from source token to the target token. (v) *MU-ST* (Zhang et al., 2022) learns an adaptive segmentation policy to detect meaningful units, which makes read/write decisions. (vi) *Baseline* is our offline-trained ST model (**B** for abbreviation). For fair comparisons, it has the same structure as MoSST.

### 5.2 Main Results

We presents the main results in Figure 4<sup>6</sup>. Compared with the online models SimulSpeech, RealTranS, and ITST, our offline model (baseline) achieves higher translation quality with high latency as it encodes bidirectional context information during training, however, in the low latency region, it performs poorly due to the input mismatch between offline-training and online-decoding.

**B + FAI** With the ability to reduce this mismatch,

<sup>3</sup><https://github.com/pytorch/fairseq>

<sup>4</sup><https://github.com/mjpost/sacrebleu>

<sup>5</sup><https://github.com/facebookresearch/SimulEval>

<sup>6</sup>The extended results for other latency metrics (AP and DAL) are described in Appendix D.5.

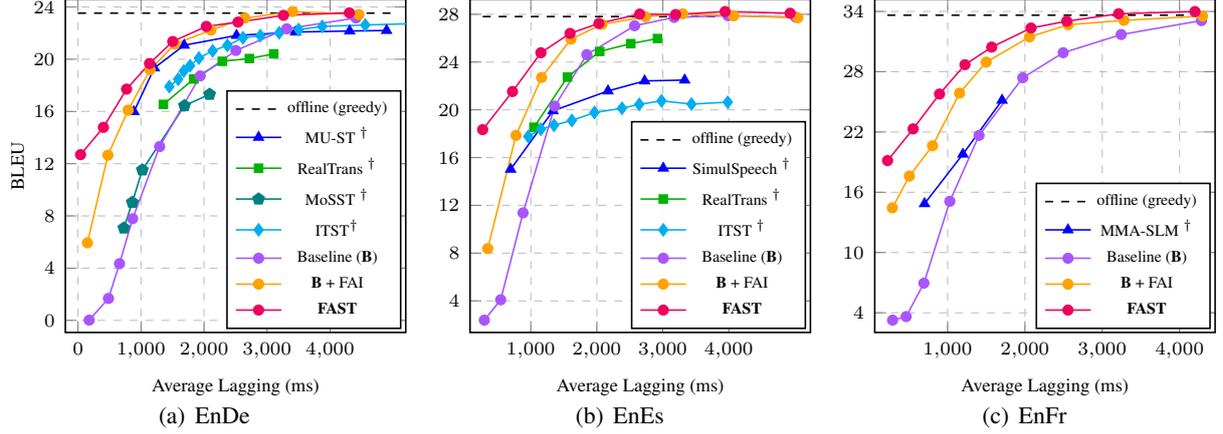


Figure 4: The translation quality (BLEU) against the latency metrics (AL) on the tst-COMMON set of MuST-C EnDe, EnEs, and EnFr dataset.  $\dagger$  denotes that the results are obtained from corresponding papers. offline is the offline performance of teacher model (offline-trained ST) by greedy search. The curve corresponding to **B** is the online performance of the teacher model using vanilla *wait-k* policy. The curve corresponding to **B + FAI** is the online performance of the teacher model with our FAI strategy. The curve corresponding to **FAST** is the online performance of our student model with the FAI strategy, *i.e.*, FAD + FAI.

FAI is directly applied for our offline (baseline) model and can achieve higher BLEU in all latency regions. In particular, it outperforms our most compatible baseline **B** by large margins in lower latency regions (when AL is less than 1000ms), with improvements over 6 BLEU in both EnDe and EnEs, 10 BLEU in EnFr.

**FAST** (FAD + FAI) Furthermore, our FAST achieves the best trade-off between translation quality and latency, especially at extremely low latency region (AL is about 200ms,  $k = 1$ ), achieving the improvements of 6 BLEU in EnDe, 10 BLEU in EnEs, and 4 BLEU in EnFr compared to **B + FAI**. It indicates that FAST can effectively mitigate the input mismatch between offline-training and online-decoding. In addition, our method achieves comparable translation quality with full-speech translation at middle latency (at AL around 2000ms), especially for EnEs.

### 5.3 Ablation Study

In this section, we study the effectiveness of our methods. All ablation results are obtained from the MuST-C EnDe tst-COMMON set. The results are shown in Figure 5.

(1) *w/o*  $\mathcal{L}_{KD}^{W2V2}$ : If removing the  $\mathcal{L}_{KD}^{W2V2}$ , the translation quality drops by 1-2 BLEU in all latency regions, including high latency region. This demonstrates optimizing  $\mathcal{L}_{KD}^{W2V2}$  can guarantee the full speech translation.

(2) *w/o*  $\mathcal{L}_{KD}^{CIF}$ : If removing the  $\mathcal{L}_{KD}^{CIF}$ , the transla-

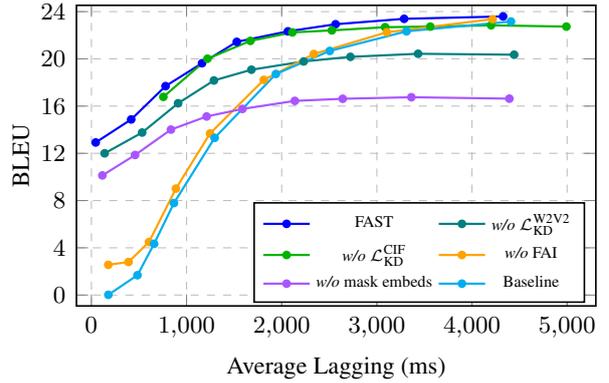


Figure 5: Ablation study of our method on the tst-COMMON set of MuST-C EnDe dataset. The observed points in the plots represent *wait-k* policy with  $k = \{1, 3, 5, 7, 9, 12, 15, 20, 30\}$ .

tion quality will be slightly degraded. However, we observe that the distances between two consecutive acoustic boundaries become larger. For example, the AL of this variant at *wait-1* is greater than 750, but the AL of the other variants at *wait-1* is approximately 150. As expected, optimizing  $\mathcal{L}_{KD}^{CIF}$  can ensure the correct acoustic boundaries.

(3) *w/o FAI*: In this variant, we use the student model by FAD with vanilla *wait-k* policy for streaming inference (*i.e.*, inference without mask tokens). However, FAD training considers mask tokens as student input, so this mismatch leads to significant performance degradation in low and middle latency regions. This indicates that our FAD and FAI should be used together to achieve better

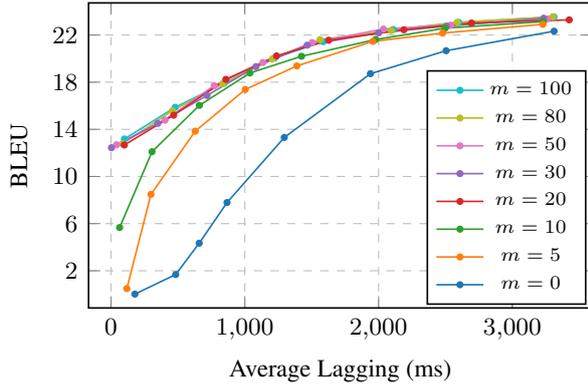


Figure 6: Effect on BLEU-AL curve of FAST w.r.t.  $m$ .

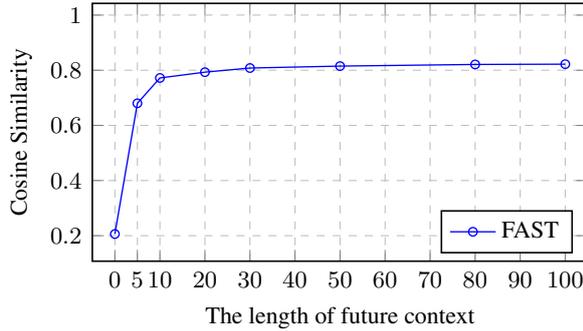


Figure 7: Effect on the  $\bar{s}_{-1}$  w.r.t.  $m$ .

streaming performance.

(4) *w/o mask embeddings*: During training and inference, our model appends  $m$  mask tokens into streaming speech tokens as the pseudo future contexts. In this variant, we remove the mask tokens during both training and inference. Even though no mismatch, we still observe a significant drop in translation quality, especially for high latency. This result indicates that the pseudo future contexts can enhance the streaming speech representations.

#### 5.4 How much future context is needed?

To answer this question, we explore the FAST (FAD + FAI) with different lengths of future context. Figure 6 shows the overall results.  $m = 0$  means the offline system without distillation. The offline system inherits the mismatch problem, but our method gradually improves the performance as  $m$  increasing from 0 to 20. Since we found only the representation of last 10 positions is poor (in Section 3), FAST obtains similar BLEU-AL curve when  $m$  is significantly larger than 10, *e.g.*, 20-100.

After the FAD training, we investigate the representation of the last position (before mask tokens) by  $\bar{s}_{-1}$  in Eq. (2) w.r.t.  $m$ . The results are shown in Figure 7. We observe that 1) as  $m$  increases, the

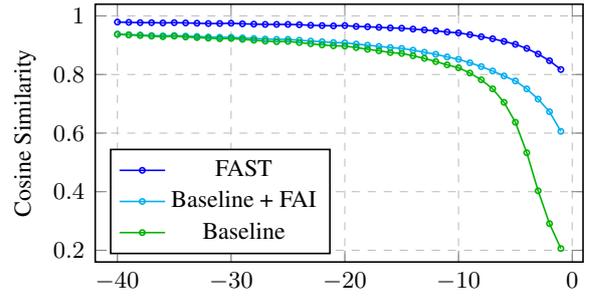


Figure 8: Effect on the average cosine similarity  $\bar{s}_{-t'}$  of the streaming speech representations at the end positions (before mask tokens). After applying FAI and FAST, the representations of the end positions are improved.

streaming speech representation of the last position becomes better; 2) the curves of the cosine similarity becomes flattened when  $m > 10$  significantly. This is consistent with the trend in Figure 6.

#### 5.5 Analysis on The Representation Gap

Figure 8 plots the changes of average cosine similarity  $\bar{s}_{-t'}$  in Eq. (2) of the last 40 positions (before mask tokens) in the streaming speech after applying the FAI or FAST (FAD + FAI). They achieve at least 0.6 and 0.8 cosine similarity at the last position, respectively. The baseline only has the  $< 0.6$  cosine similarity for the last 4 positions and only 0.2 for the last position. It indicates that the representations with FAI are closer to those of the full speech, especially at the ending positions, and FAD training can further close this gap.

## 6 Conclusion

In this paper, we examine streaming speech translation from a new perspective. We investigate the effects of the input mismatch between offline-training and online-decoding. We find that the representations at the ending positions in the streaming input are particularly poor, directly impacting the translation quality. We propose FAST, which introduces future contexts to improve these representations during training and testing via FAD and FAI, respectively. Experiments and analysis demonstrate their effectiveness in bridging the representation gap between full speech encoding and partial streaming encoding. Furthermore, our methods can be generally beneficial to streaming speech translation models that are based on Wav2Vec2.0. In the future, we will explore the relevant method independent on Wav2Vec2.0.

## 7 Limitations

Our proposed method is built upon the Wav2Vec2.0 model, whose superior representation power has been shown to enhance the performance of offline ST models. Nevertheless, it should be noted that the parameters of Wav2Vec2.0 model are considerably large, approximately 95M. As a result, this may lead to increased computational costs during training and inference. As the future work mentioned in our conclusion, we will explore the relevant method independent on Wav2Vec2.0.

The CIF module for detecting the acoustic boundary is optimized from the weakly supervised signal – total length of text tokens. In streaming inference, the boundary detector is not guaranteed to predict accurate boundaries. In other words, it is not guaranteed to align each text token with detected boundaries during online inference. However, due to the good performance of overall translation quality, we hypothesize that these boundaries may represent some meaningful acoustic units. It should be another future work to explore the underlying meaning.

## References

- Junyi Ao, Rui Wang, Long Zhou, Chengyi Wang, Shuo Ren, Yu Wu, Shujie Liu, Tom Ko, Qing Li, Yu Zhang, Zhihua Wei, Yao Qian, Jinyu Li, and Furu Wei. 2022. [SpeechT5: Unified-modal encoder-decoder pre-training for spoken language processing](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5723–5738, Dublin, Ireland. Association for Computational Linguistics.
- Naveen Arivazhagan, Colin Cherry, Wolfgang Macherey, Chung-Cheng Chiu, Semih Yavuz, Ruoming Pang, Wei Li, and Colin Raffel. 2019. [Monotonic infinite lookback attention for simultaneous machine translation](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1313–1323, Florence, Italy. Association for Computational Linguistics.
- Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. [wav2vec 2.0: A framework for self-supervised learning of speech representations](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 12449–12460. Curran Associates, Inc.
- Chih-Chiang Chang and Hung yi Lee. 2022. [Exploring Continuous Integrate-and-Fire for Adaptive Simultaneous Speech Translation](#). In *Proc. Interspeech 2022*, pages 5175–5179.

- 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*, pages 4618–4624, Online. Association for Computational Linguistics.
- Colin Cherry and George Foster. 2019. Thinking slow about latency evaluation for simultaneous machine translation. *arXiv preprint arXiv:1906.00048*.
- Kyunghyun Cho and Masha Esipova. 2016. Can neural machine translation do simultaneous translation? *arXiv preprint arXiv:1606.02012*.
- Hexuan Deng, Liang Ding, Xuebo Liu, Meishan Zhang, Dacheng Tao, and Min Zhang. 2023. Improving simultaneous machine translation with monolingual data. In *Proceedings of AAAI*.
- Mattia A. Di Gangi, Roldano Cattoni, Luisa Bentivogli, Matteo Negri, and Marco Turchi. 2019a. [MuST-C: a Multilingual Speech Translation Corpus](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2012–2017, Minneapolis, Minnesota. Association for Computational Linguistics.
- Mattia A. Di Gangi, Matteo Negri, Viet Nhat Nguyen, Amirhossein Tebbifakhr, and Marco Turchi. 2019b. [Data augmentation for end-to-end speech translation: FBK@IWSLT ‘19](#). In *Proceedings of the 16th International Conference on Spoken Language Translation*, Hong Kong. Association for Computational Linguistics.
- Lin hao Dong and Bo Xu. 2020. [Cif: Continuous integrate-and-fire for end-to-end speech recognition](#). In *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6079–6083.
- 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 (Volume 1: Long Papers)*, pages 680–694, Dublin, Ireland. Association for Computational Linguistics.
- Qianqian Dong, Mingxuan Wang, Hao Zhou, Shuang Xu, Bo Xu, and Lei Li. 2021a. Consecutive decoding for speech-to-text translation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 12738–12748.
- Qianqian Dong, Rong Ye, Mingxuan Wang, Hao Zhou, Shuang Xu, Bo Xu, and Lei Li. 2021b. Listen, understand and translate: Triple supervision decouples end-to-end speech-to-text translation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 12749–12759.





- 949 Shaolei Zhang and Yang Feng. 2022a. [Gaussian multi-](#)  
950 [head attention for simultaneous machine translation.](#)  
951 In *Findings of the Association for Computational*  
952 *Linguistics: ACL 2022*, pages 3019–3030, Dublin,  
953 Ireland. Association for Computational Linguistics.
- 954 Shaolei Zhang and Yang Feng. 2022b. [Information-](#)  
955 [transport-based policy for simultaneous translation.](#)  
956 In *Proceedings of the 2022 Conference on Empirical*  
957 *Methods in Natural Language Processing*, pages 992–  
958 1013, Abu Dhabi, United Arab Emirates. Association  
959 for Computational Linguistics.
- 960 Shaolei Zhang and Yang Feng. 2022c. [Modeling dual](#)  
961 [read/write paths for simultaneous machine transla-](#)  
962 [tion.](#) In *Proceedings of the 60th Annual Meeting of*  
963 *the Association for Computational Linguistics (Vol-*  
964 *ume 1: Long Papers)*, pages 2461–2477, Dublin,  
965 Ireland. Association for Computational Linguistics.
- 966 Shaolei Zhang and Yang Feng. 2022d. [Reducing posi-](#)  
967 [tion bias in simultaneous machine translation with](#)  
968 [length-aware framework.](#) In *Proceedings of the 60th*  
969 *Annual Meeting of the Association for Computational*  
970 *Linguistics (Volume 1: Long Papers)*, pages 6775–  
971 6788, Dublin, Ireland. Association for Computational  
972 Linguistics.
- 973 Shaolei Zhang, Yang Feng, and Liangyou Li. 2021.  
974 Future-guided incremental transformer for simulta-  
975 neous translation. In *Proceedings of the AAAI Con-*  
976 *ference on Artificial Intelligence*, volume 35, pages  
977 14428–14436.
- 978 Jiawei Zhao, Wei Luo, Boxing Chen, and Andrew  
979 Gilman. 2021. [Mutual-learning improves end-to-](#)  
980 [end speech translation.](#) In *Proceedings of the 2021*  
981 *Conference on Empirical Methods in Natural Lan-*  
982 *guage Processing*, pages 3989–3994, Online and  
983 Punta Cana, Dominican Republic. Association for  
984 Computational Linguistics.
- 985 Qinpei Zhu, Renshou Wu, Guangfeng Liu, Xinyu Zhu,  
986 Xingyu Chen, Yang Zhou, Qingliang Miao, Rui  
987 Wang, and Kai Yu. 2022. [The AISP-SJTU simultane-](#)  
988 [ous translation system for IWSLT 2022.](#) In *Proceed-*  
989 *ings of the 19th International Conference on Spoken*  
990 *Language Translation (IWSLT 2022)*, pages 208–215,  
991 Dublin, Ireland (in-person and online). Association  
992 for Computational Linguistics.

---

**Algorithm 1** Pseudocode of FAI strategy in a PyTorch-like style.

---

```
# model: an offline-trained ST model consists of a acoustic encoder Wav2vec2.0, a token boundary detector
, a semantic encoder, and a decoder
# m: mask length, K: wait lagging, audio: audio waveform
# mask_emb: pre-trained mask embedding in Wav2vec

N = 0 # the number of source text tokens
x = [] # streaming audio prefix
y = [] # translations
mask_embs = mask_emb.repeate(m, 1) # mask embeddings:  $m \times d$ 
while y[-1] != "<eos>":
    if x == audio: # audio has been read
        y = y + model(a,y) # write new target token
    elif N - len(y) < K: # wait K detected source tokens
        x = x + read(audio) # incrementally read audio
        c = model.wav2vec2.cnn(x) # audio tokens  $\tau \times d$ 

        c = torch.cat((c, mask_embs), dim=0) # concatenate audio tokens and mask embeddings,  $(\tau + m) \times d$ 
        a = model.wav2vec2.encoder(c) # audio representations,  $(\tau + m) \times d$ 
        a = a[:a.shape[0] - m,:] # discard the predicted representations,  $\tau \times d$ 

        if model.token_detector(a): # source text token boundary is detected
            N += 1
    else:
        h = model.semantic_encoder(a)
        y = y + model.decoder(h, y) # write new target token
```

---

## A Data Statistics

We evaluate our model on MuST-C V1 English-German (EnDe), English-Spanish (EnEs) and English-French (EnFr) datasets (Di Gangi et al., 2019a). For training set, we follow Dong et al. (2022) to filter out short speech of less than 1000 frames (62.5ms) and long speech of more than 480,000 frames (30s). The statistics of different language pairs are illustrated in Table 1.

split	EnDe	EnEs	EnFr
train	225,271	260,041	269,248
dev	1,418	1,312	1,408
tst-COMMON	2,641	2,502	2,632

Table 1: Number of samples for each split of MuST-C datasets.

## B Additional Preliminary Analysis

### B.1 Which part of streaming speech representation is worse?

To further verify that only the representation of the end position in streaming speech is poor, we calculate the cosine similarity  $s_{t,t'}$  between the speech representation at the  $t'$ -th ( $t' \leq t$ ) position in the  $t$ -th streaming audio input  $\hat{\mathbf{x}}_t$  and the speech representation at the same position in the full encoding. Then we average the cosine similarities over the sentences in dataset  $\mathcal{B}$  to obtain robust statistics.

$$\text{For } t' \leq t, \quad \bar{s}_{t,t'} = \frac{1}{|\mathcal{B}_t|} \sum_{\mathbf{x} \in \mathcal{B}_t} s_{t,t'}(\mathbf{x}) = \frac{1}{|\mathcal{B}_t|} \sum_{\mathbf{x} \in \mathcal{B}_t} \cos(\hat{a}_{t,t'}, a_{t'}), \quad (12)$$

where  $\mathcal{B}_t = \{\mathbf{x} : |\mathbf{x}| \geq t\}$  contains the audio inputs with length no shorter than  $t$ .

We empirically compare the averaged cosine similarity at the beginning, middle, and end positions of the speech representations. Figure 9 shows  $\bar{s}_{t,t'}$  of the first three ( $t' = 1, 2, 3$ ), middle three ( $t' = \lfloor \frac{1+t}{2} \rfloor - 1, \lfloor \frac{1+t}{2} \rfloor, \lfloor \frac{1+t}{2} \rfloor + 1$ ), and last three ( $t' = t - 2, t - 1, t$ ) positions for each encoding step  $t$ . At the beginning and middle positions, the averaged cosine similarity  $\bar{s}_{t,t'}$  is greater than 0.8 except  $t' = 1$ , indicating that the representations at such positions in the partial streaming input are close to those in the full speech. Note that  $t' = 1$  with a slightly lower similarity won't hurt the performance much, because in practice it is almost impossible to apply *wait-1* policy (only read 20ms speech input) in streaming ST. However, the  $\bar{s}_{t,t'}$  declines significantly for the end positions, especially for the last one. In addition, we

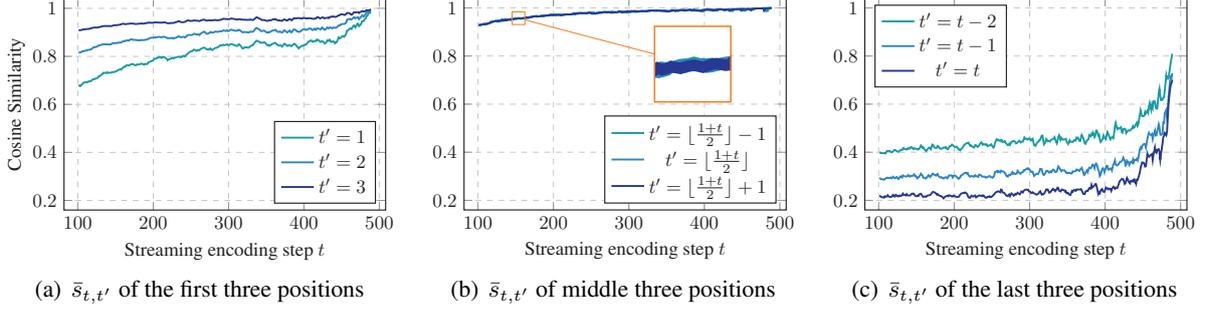


Figure 9: The average cosine similarity  $\bar{s}_{t,t'}$  of the first three ( $t' = 1, 2, 3$ ), middle three ( $t' = \lfloor \frac{1+t}{2} \rfloor - 1, \lfloor \frac{1+t}{2} \rfloor, \lfloor \frac{1+t}{2} \rfloor + 1$ ), and last three ( $t' = t - 2, t - 1, t$ ) positions for each encoding step  $t$ .

observe that as  $t$  becomes larger, the streaming input will gradually approximate the full speech input, then the gap of the speech representation between the offline and the online input becomes smaller. We conclude that **the representations of the end position in the streaming speech are particularly inferior.**

## B.2 Does the poor representation at the last positions of streaming speech affect streaming ST performance?

To answer this question, we only calculate the average cosine similarity in the last position for each sample.

$$\forall \mathbf{x}, \quad \bar{s}_{-1}(\mathbf{x}) = \frac{1}{T} \sum_{t=1}^{t=T} \cos(\hat{a}_{t,t}, a_t), \quad (13)$$

$\bar{s}_{-1}(\mathbf{x})$  reflects the degree of deterioration of the representation at the last position of the streaming speech. We sort the dataset by the value of the degree and divide them evenly into 5 groups to ensure enough samples in each group. The translation quality of each group is shown in Figure 10. The performance of streaming ST drops close to 10 points as the representation at the last position of the streaming speech becomes worse, while the full-sentence ST fluctuates less than 4 points. In addition, the performance gap between the streaming ST and the full-sentence ST becomes larger as the representation at the last position gets worse. In the worse group, the streaming ST is 12.41 points lower than the full-sentence ST. Therefore, we conclude that **the poor representation at the end position of the streaming speech has a strong effect on the translation quality.**

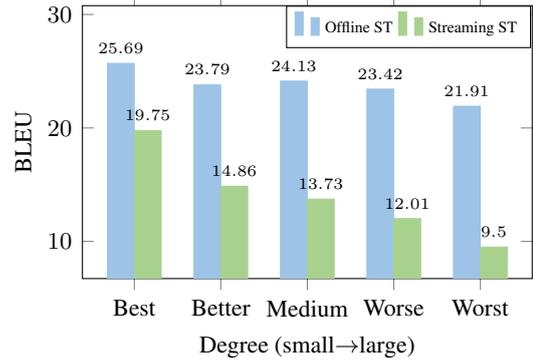


Figure 10: Performance with degree of deterioration of the representation at the last position of the streaming speech.

## C Details of Offline Training

We use an Adam optimizer with learning rate  $1e^{-4}$  and warmup step  $10k$ . We decay the learning rate with inverse square root schedule.

The offline ST model is first trained by a multi-task learning, including ASR and ST tasks. A language identity tag is prepended to the target sentence for indicating which task is learned. In this stage, the CIF module which is used to detect the acoustic boundary is deactivated, in other words, the CIF module is not trained. The main purpose is to learn a better decoder, i.e., a well-trained language model. Then, we activate the CIF module such that its parameters are trainable, and continue to train for another several epochs. In this stage, only the ST task is learned.

Monotonic Level	Easy	Medium	Hard	AL
Offline (greedy)	26.38	23.22	21.26	-
Baseline	18.88	12.95	10.38	1295
+ FAI	23.88 <sup>+5.00</sup>	18.99 <sup>+6.04</sup>	16.45 <sup>+6.07</sup>	1143
FAST	24.44 <sup>+5.56</sup>	19.89 <sup>+6.94</sup>	16.53 <sup>+6.15</sup>	1135

Table 2: Performance (BLEU) on different monotonic levels on test set of MuST-C EnDe.

## D Additional Experiments 1046

### D.1 Why we use AL rather than $k$ ? 1047

In our presented results, we plot the BLEU *v.s.* AL rather than  $k$ . We argue that  $k$  is not a fair metric to evaluate the latency. In text streaming translation, different tokenization (*e.g.*, different number of BPE operations) will lead to different token boundaries for the same sentence. It indicates the  $k$  tokens do not necessarily represent the same partial sentence for different BPE methods. This situation becomes even severer for speech streaming translation. As we have a source text token boundary detector in our model, the first  $k$  detected text tokens will represent different lengths of audio frames for different input audios. To be precise, the wait- $k$  policy used in our streaming speech translation is actually wait- $k$  detected tokens policy. Therefore, we prefer to use AL rather than  $k$  as the latency metric in our experiments. 1051-1055

### D.2 What examples are improved by our strategies? 1056

For tst-COMMON on MuST-C EnDe, we use awesome-align<sup>7</sup> (Dou and Neubig, 2021) to identify the token-level alignment between source transcription and target translation following Zhang and Feng (2022d). First, we define the source-to-target alignment position shift as  $\max\{0, i - j\}$ , where the  $i$ th source token is aligned to the  $j$ th target token. If  $i - j$  is large, it means in order to translate the  $j$ th target token, the model may need to read more until seeing the  $i$ th source token. Then we calculate the monotonic level of each example as the averaged alignment position shift over the number of aligned tokens, *i.e.*, 1057-1063

$$\text{monotonic\_level} = \frac{1}{|\text{aligned\_pairs}|} \sum_{(i,j) \in \text{aligned\_pairs}} \max\{0, i - j\} \quad (14) \quad 1064$$

We evenly divide the test set into three groups according to different monotonic levels. For each group, we evaluate different inference methods and report the results in Table 2. As we explained in D.1, it is almost impossible to guarantee the same AL for different inference methods. For a fair comparison, we try our best to set the AL of different methods to be approximately equal. We can see our inference strategies show a significant advantage on the non-monotonic examples (medium and hard groups). 1065-1069

### D.3 How important of the Wav2Vec2.0? 1070

As we mentioned in the main text, the special audio token “mask” in Wav2Vec2.0 is pre-trained on the Librispeech dataset to reconstruct the corresponding feature conditional on unmasked context via the contrastive task. In our experiments, we didn’t include contrastive learning as the auxiliary task in the downstream ST training. And in our FAI inference, we directly leverage the mask embeddings as the future context by appending them to the streaming input. However, we found the speech representations after ST training becomes even better. Particularly, we calculate the cosine similarity between every predicted future representation and full speech representations at the same position, and the results are illustrated in Figure 11. On either the Librispeech or the MuST-C audio test set, the fine-tuned Wav2Vec2.0 can produce better speech representations from the masking inputs. 1071-1079

### D.4 Why are all predicted features discarded? 1080

<sup>7</sup><https://github.com/neulab/awesome-align>

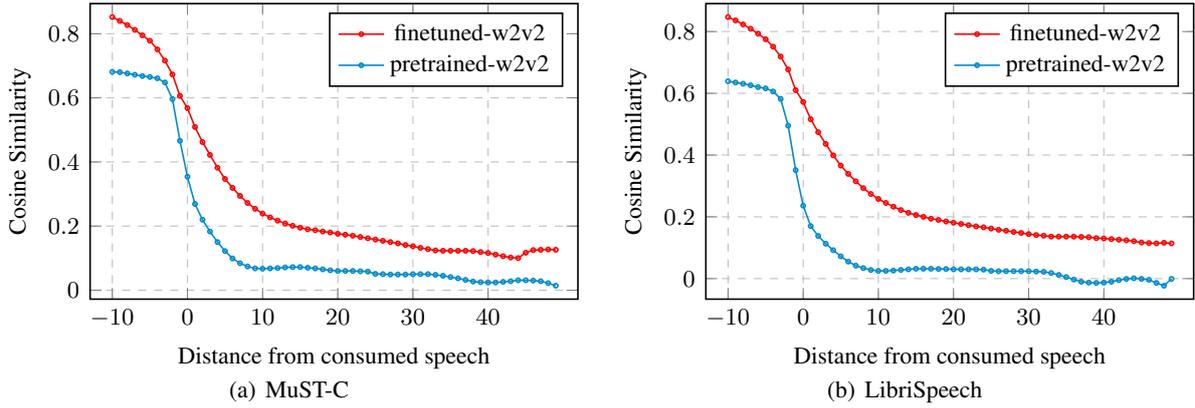


Figure 11: We measure the accuracy of predicted context by calculating the cosine similarity between every predicted future representation and full speech representations at the same position.

1081 In FAI strategy, all the output representations corre-  
 1082 sponding to the  $m = 50$  masking tokens will be dis-  
 1083 carded, because we have demonstrated that the repre-  
 1084 sentations at the ending positions are inferior. However,  
 1085 as shown in 11, the first 10 predicted representations are  
 1086 not as bad as the next 40. Therefore, on the EnDE test  
 1087 set, we also conduct another streaming ST inference  
 1088 by appending different numbers of predicted context  
 1089 to the original speech representations. We use discard  
 1090 rate  $p$  to measure the number of appending features.  
 1091 When  $p = 1.0$ , all predicted features are discarded and  
 1092 it reduces to the standard FAI inference. In Figure 12,  
 1093 we compare the streaming speech translation quality  
 1094 between regular FAI and its variant. It is concluded that  
 1095 the predicted future context is too noisy and harmful to the performance.

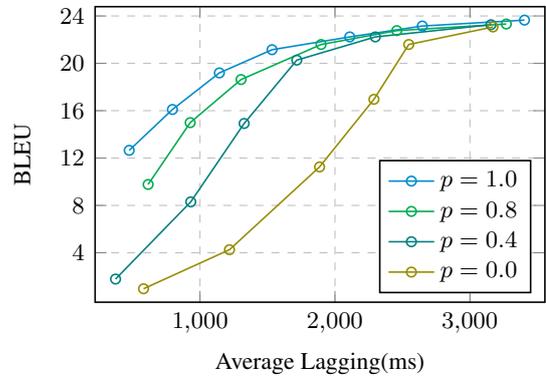


Figure 12: BLEU v.s. AL on different  $p$ .

## 1096 D.5 Additional Results on EnDe/Es and EnFr

1097 In this section, we evaluate our methods with other latency metrics AP and DAL. The AP-BLEU and  
 1098 DAL-BLEU curves on the MuST-C EnDe, EnEs, and EnFr tst-COMMON sets are shown in Figure 13.  
 1099 For three language pairs, our proposed methods can consistently improve the baseline by a large margin.

## 1100 E Numeric Results for the Figures

1101 We also provide the numeric results for Figures 4 and 13 in Tables 3, and for Figures 5 in Table 4, and for  
 1102 Figures 4 in Table 5, for Figure 6 in Table 6.

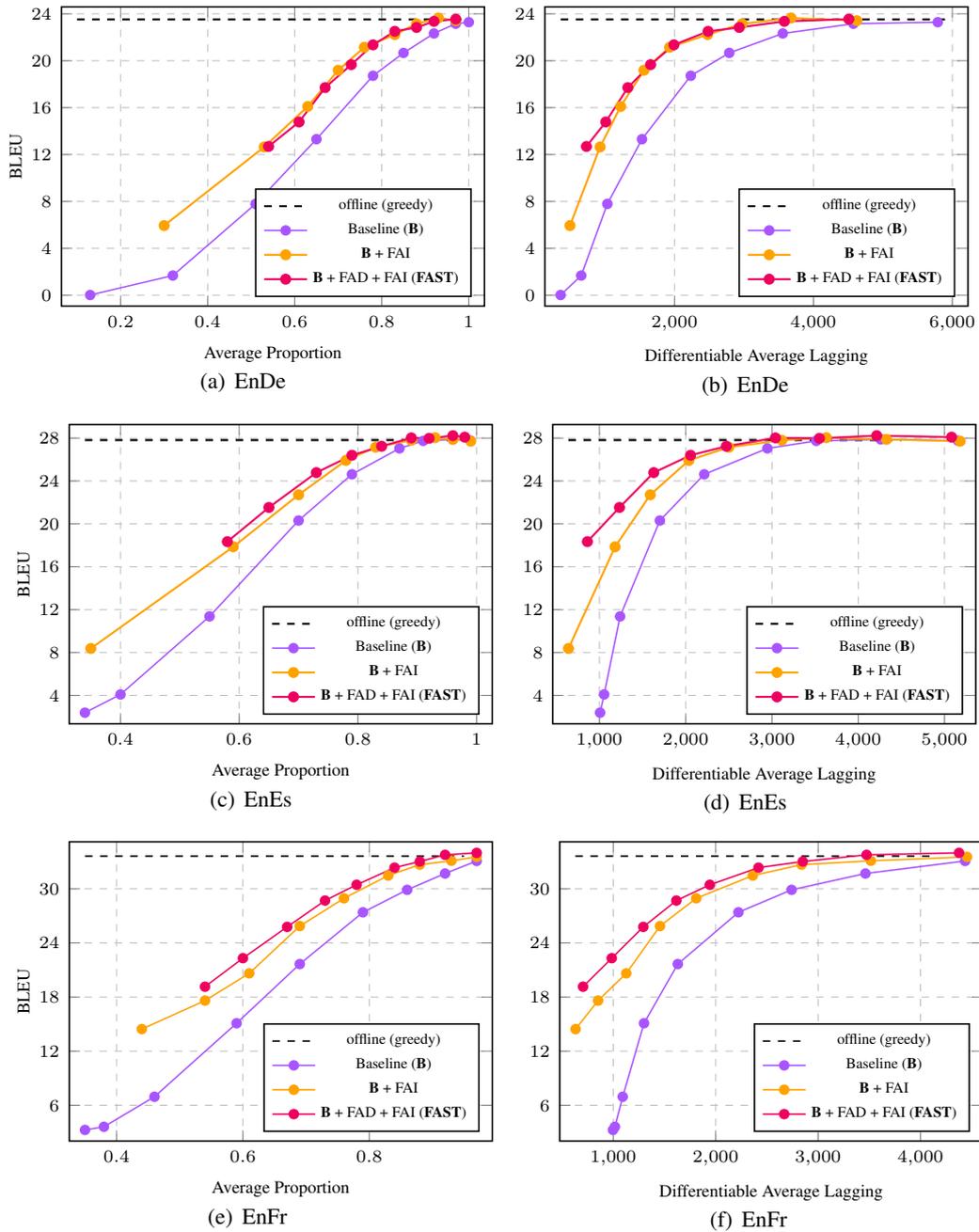


Figure 13: The translation quality (BLEU) against the latency metrics (AP, DAL) on the tst-COMMON set of MuST-C EnDe, EnEs and EnFr dataset.

Model	Lagging ( $k$ )	En-De				En-Es				En-Fr			
		AL	AP	DAL	BLEU	AL	AP	DAL	BLEU	AL	AP	DAL	BLEU
Baseline	1	178	0.13	359	0.02	295	0.34	1007	2.39	288	0.35	997	3.27
	3	483	0.32	656	1.68	543	0.4	1054	4.09	463	0.38	1016	3.62
	5	659	0.42	821	4.34	882	0.55	1239	11.37	693	0.46	1092	6.95
	7	867	0.51	1032	7.79	1361	0.7	1700	20.31	1028	0.59	1300	15.1
	9	1295	0.65	1531	13.31	1848	0.79	2215	24.62	1406	0.69	1630	21.66
	12	1939	0.78	2234	18.72	2572	0.87	2947	27.04	1972	0.79	2222	27.4
	15	2505	0.85	2788	20.67	3171	0.91	3513	27.74	2495	0.86	2741	29.89
	20	3312	0.92	3559	22.33	3988	0.96	4260	27.88	3245	0.92	3462	31.7
	30	4410	0.97	4576	23.16	5012	0.99	5157	27.76	4283	0.97	4435	33.09
+ FAI	1	150	0.3	494	5.94	347	0.35	641	8.38	285	0.44	632	14.45
	3	475	0.53	928	12.65	775	0.59	1181	17.86	505	0.54	852	17.61
	5	796	0.63	1223	16.1	1162	0.7	1589	22.71	805	0.61	1127	20.63
	7	1143	0.7	1559	19.19	1608	0.78	2037	25.92	1154	0.69	1456	25.87
	9	1534	0.76	1928	21.15	2076	0.83	2500	27.15	1498	0.76	1810	28.95
	12	2109	0.83	2476	22.23	2736	0.89	3114	27.8	2060	0.83	2362	31.47
	15	2647	0.88	2974	23.15	3301	0.93	3630	28.04	2559	0.88	2838	32.68
	20	3404	0.93	3678	23.65	4072	0.96	4328	27.88	3280	0.93	3515	33.11
	30	4457	0.97	4625	23.42	5045	0.99	5181	27.71	4297	0.97	4454	33.54
FAST	1	41	0.54	731	12.69	270	0.58	860	18.34	223	0.54	705	19.15
	3	403	0.61	1009	14.78	722	0.65	1232	21.53	554	0.6	985	22.31
	5	771	0.67	1327	17.71	1152	0.73	1629	24.78	895	0.67	1293	25.78
	7	1135	0.73	1655	19.67	1594	0.79	2056	26.4	1224	0.73	1616	28.7
	9	1503	0.78	1991	21.36	2031	0.84	2471	27.24	1570	0.78	1943	30.45
	12	2036	0.83	2483	22.51	2650	0.89	3040	28.02	2079	0.84	2418	32.35
	15	2539	0.88	2932	22.84	3194	0.92	3550	27.98	2541	0.88	2850	33.03
	20	3260	0.92	3581	23.36	3943	0.96	4214	28.23	3212	0.92	3473	33.77
	30	4305	0.97	4510	23.55	4928	0.98	5082	28.09	4199	0.97	4376	33.99

Table 3: Numeric results on MuST-C EnDe, EnEs, and EnFr tst-COMMON set (Figure 4 and 13).

Lagging ( $k$ )	w/o $\mathcal{L}_{KD}^{W2V2}$		w/o $\mathcal{L}_{KD}^{CIF}$		w/o FAI		w/o mask embeds	
	AL	BLEU	AL	BLEU	AL	BLEU	AL	BLEU
1	139	12	756	16.78	177	2.56	115	10.13
3	533	13.76	1220	20.01	390	2.8	459	11.86
5	911	16.24	1671	21.52	605	4.49	836	14.01
7	1288	18.17	2112	22.24	888	9.01	1211	15.12
9	1682	19.08	2527	22.41	1247	13.68	1588	15.76
12	2231	19.78	3087	22.68	1812	18.22	2138	16.44
15	2722	20.17	3562	22.73	2338	20.41	2641	16.62
20	3434	20.43	4201	22.84	3105	22.25	3363	16.75
30	4443	20.35	4992	22.73	4217	23.36	4393	16.63

Table 4: Numeric results for ablation study (Figure 5).

	<b><i>MU-ST</i></b>										
	AL	1023	1424	1953	2642	3621	4453	5089	5754		
	BLEU	17.94	20.85	22.78	24.3	24.82	24.99	25.05	25.9		
	<b><i>RealTrans</i></b>										
	AL	1355	1838	2290	2720	3106					
	BLEU	16.54	18.49	19.84	20.05	20.41					
EnDe	<b><i>MoSST</i></b>										
	AL	728	862	1021	1689	2088					
	BLEU	7.07	9.04	11.52	16.44	17.31					
	<b><i>ITST</i></b>										
	AL	1449	1589	1678	1778	1919	2137	2371			
	BLEU	17.9	18.47	19.09	19.5	20.09	20.64	21.06			
	AL	2618	2893	3193	3501	3876	4557	5206			
	BLEU	21.64	21.8	22.02	22.27	22.51	22.62	22.71			
	<b><i>SimulSpeech</i></b>										
	AL	694	1336	2169	2724	3331					
	BLEU	15.02	19.92	21.58	22.42	22.49					
	<b><i>RealTrans</i></b>										
EnEs	AL	1047	1554	2043	2514	2920					
	BLEU	18.54	22.74	24.89	25.54	25.97					
	<b><i>ITST</i></b>										
	AL	960	1153	1351	1621	1964	2381	2643	2980	3434	3983
	BLEU	17.77	18.38	18.71	19.11	19.77	20.13	20.46	20.75	20.48	20.64
	<b><i>MMA-SLM</i></b>										
EnFr	AL	701	1197	1704							
	BLEU	14.86	19.79	25.16							

Table 5: Numeric results for baseline systems (Figure 4). The results of *MU-ST* are obtained from (Zhang et al., 2022). The results of *SimulSpeech* and *RealTrans* are obtained from (Zeng et al., 2021). The results of *MoSST* are obtained from (Dong et al., 2022). The results of *ITST* are obtained from (Zhang and Feng, 2022b). The results of *MMA-SLM* are obtained from (Indurthi et al., 2022).

Lagging ( $k$ )	$m = 5$		$m = 10$		$m = 20$		$m = 30$		$m = 50$		$m = 80$		$m = 100$	
	AL	BLEU	AL	BLEU	AL	BLEU	AL	BLEU	AL	BLEU	AL	BLEU	AL	BLEU
1	118	0.49	64	5.67	99	12.67	3	12.44	41	12.69	85	12.78	100	13.18
3	298	8.48	306	12.1	468	15.2	349	14.5	403	14.78	458	15.57	479	15.87
5	629	13.84	660	16.03	858	18.24	717	16.87	771	17.71	835	17.87	845	17.91
7	1003	17.38	1038	18.78	1237	20.23	1083	19.32	1135	19.67	1205	19.97	1225	20.07
9	1389	19.38	1424	20.2	1627	21.56	1466	21.14	1503	21.36	1562	21.61	1587	21.44
12	1957	21.46	1978	21.62	2189	22.45	2001	22.19	2036	22.51	2095	22.38	2109	22.47
15	2479	22.17	2497	22.58	2695	23.02	2507	22.75	2539	22.84	2588	23.08	2599	23.07
20	3228	22.91	3231	23.14	3425	23.29	3234	23.43	3260	23.36	3302	23.55	3311	23.54

Table 6: Numeric results for different lengths future context (Figure 6).