

FASST: Fast LLM-based Simultaneous Speech Translation

Anonymous ACL submission

Abstract

Simultaneous speech translation (SST) takes streaming speech input and generates text translation on the fly. Existing methods either have high latency due to recomputation of input representations, or fall behind of offline ST in translation quality. In this paper, we propose FASST, a fast large language model based method for streaming speech translation. We propose blockwise-causal speech encoding and consistency mask, so that streaming speech input can be encoded incrementally without recomputation. Furthermore, we develop a two-stage training strategy to optimize FASST for simultaneous inference. We evaluate FASST and multiple strong prior models on MuST-C dataset. Experiment results show that FASST achieves the best quality-latency trade-off. It outperforms the previous best model by an average of 1.5 BLEU under the same latency for English to Spanish translation.

1 Introduction

End-to-end simultaneous speech translation (SST) translates incomplete speech input into text in a different language (Ma et al., 2020b), which is widely used in multilingual conferences, live streaming and etc. Compared to offline ST where speech input is complete, SST needs to decide whether to continue waiting or to generate more translation after receiving new speech input. A common approach in building performant SST streaming models involves pretraining for offline translation and optional finetuning for simultaneous translation (Agarwal et al., 2023; Communication et al., 2023). The quality-latency trade-off of simultaneous streaming models thus heavily depends on its offline performance.

Large language model (LLM) have recently demonstrated its potential to be a strong backbone of offline E2E ST (Huang et al., 2023; Zhang et al., 2023b). However, LLM introduces larger computation overhead compared to regular-sized models

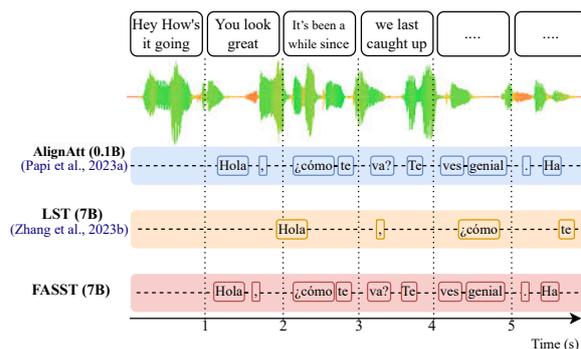


Figure 1: Simultaneous speech translation with AlignAtt-0.1B, LST-7B and our FASST-7B. The LST-7B model generates translation with significantly higher latency than AlignAtt, while our FASST-7B achieves comparable latency with it.

when applied to SST. Figure 1 shows that the computation latency of a LLM-based 7B model makes it inferior for real-time application.

The computation overhead of SST models comes from both encoding new speech input and decoding new translation. While the latter one has been heavily optimized for LLM (Pope et al., 2022; Kwon et al., 2023; Dao, 2024), the former one has not been optimized for SST. As new speech input arrives, most SST models re-encode the entire speech and start autoregressive decoding afterwards, ignoring the incremental nature of streaming speech input. More importantly, the LLM decoder needs to recompute hidden states due to the updated speech features, significantly slowing down the computation.

In this work, we propose a **FA**st LLM-based **SS**T (FASST) method to avoid recomputation while maintaining its translation quality. We develop a blockwise-causal speech encoding technique that incrementally encodes new speech input and introduce incremental LLM decoding with consistency mask. We also design a 2-stage training strategy for FASST: 1) aligning speech encoder outputs

with LLM embeddings using word-aligned contrastive loss (Ouyang et al., 2023) and 2) finetuning for SST using wait- k -stride- n policy (Zeng et al., 2021). Experiments on MuST-C dataset (Di Gangi et al., 2019) shows that our 7B model maintains competitive computation aware latency compared to 115M baselines while achieving consistent quality improvement of at least 1.5 BLEU score on English-Spanish direction.

Our contributions are:

- We propose FASST, one of the first efficient LLM-based methods for simultaneous speech translation.
- We verify FASST on MuST-C dataset and it outperforms strong prior methods by 1.5 BLEU at the same latency on English-Spanish direction.
- We further demonstrate that FASST can be generalized to other policies like hold- n and policies spending more time on encoding benefit more from FASST.

2 Related Works

End-to-End SST translates partial speech input into text in another language without generating intermediate transcription. A variety of speech segmentation techniques and policies have been proposed to optimize the quality-latency trade-off. Ren et al. (2020); Dong et al. (2022); Zeng et al. (2023); Zhang et al. (2023c) learn to segment streaming speech input by word boundaries. Zhang and Feng (2023) further learns to segment speech at moments that are beneficial to the translation. On the policy side, Ma et al. (2020b) adapts wait- k and monotonic multihead attention (MMA) from simultaneous text translation to SST model. Ma et al. (2023) further improves the numerical stability of MMA. Papi et al. (2023b) constructs source-target alignment with attention information to guide the simultaneous inference. Zhang and Feng (2022) decides whether to translate based on accumulated information of source speech. Polák et al. (2023) conducts blockwise beam search when doing incremental decoding. The translation quality of SST models depend on not only their policies, but also their offline performance (Agarwal et al., 2023). Recently LLM has been shown as a strong backbone of offline ST (Zhang et al., 2023b; Huang et al., 2023), but its computation overhead prevents it from being used in SST scenarios. FASST is one of the first LLM-based SST models with a reasonable quality-latency trade-off.

Efficient ST To reduce the computation cost of ST models, Wu et al. (2020); Ma et al. (2020c); Raffel and Chen (2023); Raffel et al. (2023) use segments and explicit or implicit memory banks to calculate self-attention only within the segment. Zhang and Feng (2023); Chen et al. (2021); Wu et al. (2021) adopt unidirectional attention during speech encoding. These methods focus on encoder-side optimization and can be integrated with FASST.

Translation with LLM While LLMs are capable of zero-shot machine translation (Brown et al., 2020; OpenAI, 2023; Touvron et al., 2023a,b), their performance can be further improved via in-context learning (Vilar et al., 2023; Zhang et al., 2023a), supervised and semi-supervised finetuning (Rothe et al., 2020; Yang et al., 2023; Zhang et al., 2023d; Xu et al., 2023). For simultaneous machine translation (SMT), Guo et al. (2024) propose a collaborative translation model with two LLM agents and Koshkin et al. (2024) design a finetuning strategy by adding a special "wait" token. Raffel et al. (2024) propose SimulMask to mask token connections under certain policy. SimulMask is a concurrent work with us and only works on text translation.

3 The FASST Method

In this section, we first review the problem formulation of simultaneous speech translation (SST) and then describe the architecture of our proposed model, FASST, followed by its training and inference strategies.

3.1 Problem Formulation

Simultaneous speech translation (SST) needs to generate translations while receiving streaming speech input. Let $S = (s_1, s_2, \dots, s_{|S|})$ be a speech waveform where s_i are real numbers. The streaming speech input is cut into segments S_1, S_2, \dots and the SST model P_θ needs to emit partial translations T_1, T_2, \dots after receiving each of them,

$$T_i \sim P_\theta(\cdot \mid S_{\leq i}, T_{< i}). \quad (1)$$

T_i can be an empty string, indicating that the SST model needs more speech input to continue the translation. After receiving all inputs S_1, S_2, \dots, S_m and emitting all translations T_1, T_2, \dots, T_m , we obtain the final translation $T = \bigoplus_{i=1}^m T_i$ by concatenating all partial ones.

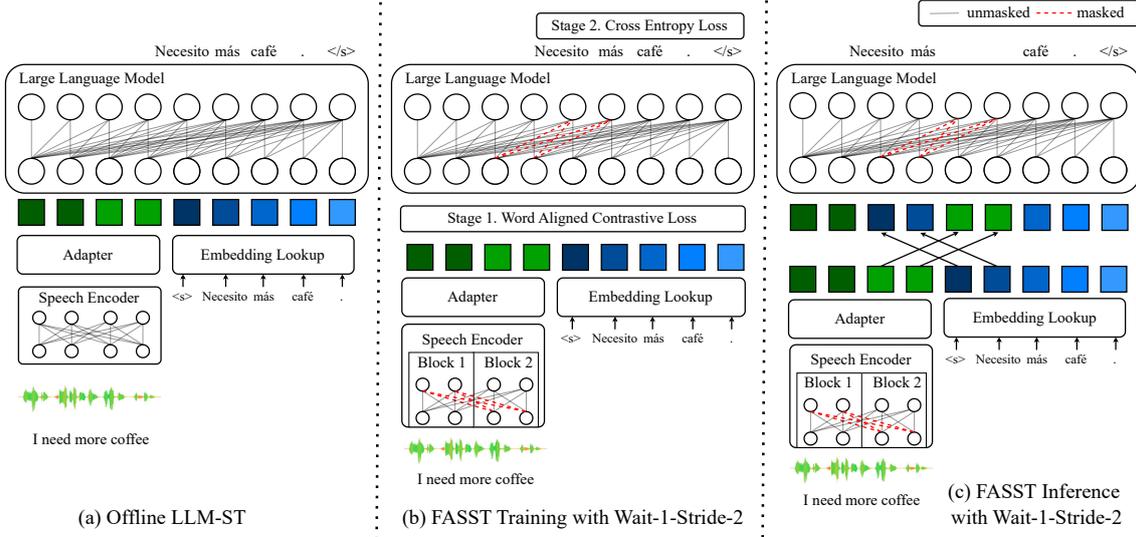


Figure 2: Overview of FASST. (a) shows the offline translation of LLM-based ST model. (b) depicts the 2-stage training pipeline of FASST. Stage 1 aligns adapter output with LLM embedding and stage 2 finetunes for simultaneous translation using wait- k -stride- n policy. (c) illustrates the simultaneous inference procedure of FASST with incremental speech encoding and LLM decoding with consistency mask.

The objective of SST is to emit high-quality translation at a low latency. Quality is evaluated by comparing generated T with the ground-truth T^* , while latency is evaluated based on the amount of lagging of each generated word. In this paper, we consider the computation-aware latency of SST models.

3.2 Model Architecture

As shown in Figure 2, our model is composed of a speech encoder, an adapter and a LLM decoder.

Blockwise-Causal Speech Encoder (BCSE) extracts contextualized acoustic features from the raw waveform incrementally. It consists of several casual convolutional layers as the audio feature extractor and a blockwise-causal Transformer Encoder as the contextual encoder.

Our causal convolutional layers are built upon non-causal ones. Denote $H_{in} \in \mathbb{R}^{l \times d}$ as the input vectors to non-causal convolution $\text{Conv}(\cdot)$ with kernel size w . We add additional zero padding $\text{Pad} \in \mathbb{R}^{(w/2-1) \times d}$ to its left so that each output vector only depends on input vectors to its left, and remove the last $w/2 - 1$ states to keep its output length the same as before,

$$H_{out} = \text{Conv}(\text{Pad} \oplus H_{in})_{:-w/2+1}. \quad (2)$$

Besides, we apply blockwise-causal masking to Transformer Encoder. Define attention mask M of

speech encoder as follows

$$M_{j_Q, j_K} = \begin{cases} 0 & \lfloor \frac{j_Q}{b} \rfloor \geq \lfloor \frac{j_K}{b} \rfloor \\ -\infty & \text{otherwise} \end{cases} \quad (3)$$

where b is the block size, i.e., the number of hidden states of the speech encoder corresponding to one segment, and j_Q, j_K are row indices of query matrix Q and key matrix K . The attention output of speech encoder during training can then be written as

$$O = \text{Softmax} \left(\frac{QK^T}{\sqrt{d}} + M \right) V, \quad (4)$$

where V is the value matrix.

Adapter receives speech encoder outputs and converts them to the LLM embedding space. It consists of two causal convolutional layers to reduce the length of speech encoder outputs by four and one linear layer to project features into the LLM embedding space. We call the adapter outputs as speech embeddings,

$$E_{\leq i}^s = \text{Adapter}(\text{BCSE}(S_{\leq i})). \quad (5)$$

LLM receives speech embeddings and embeddings of previously generated tokens to decode autoregressively according to a wait- k -stride- n policy π .

$$T_i \sim \text{LLM}(\cdot | E_{\leq i}^s, T_{< i}, \pi). \quad (6)$$

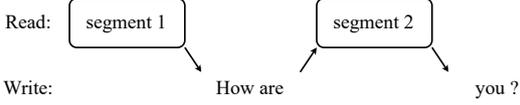


Figure 3: Example of wait-1-stride-2. It waits for 1 segment at the beginning and then alternate between generate 2 words (including punctuations) and reading new segment.

Wait- k -stride- n policy waits for k speech segments at the beginning and then alternate between generating n words and reading new segment. Figure 3 shows an example of wait-1-stride-2.

3.3 Training

As shown in Figure 2 (b), we employ a 2-stage approach to train our model.

Stage 1. Speech-text alignment. We align the speech embedding with LLM input embedding using word-aligned contrastive (WACO) loss. Both transcription embeddings E^t and speech embeddings E^s are grouped into word embeddings W^t and W^s by word boundaries. Word boundaries of speech are obtained through Montreal Forced Aligner¹. We treat speech and transcription embeddings of the same word as positive pair and others as negative pairs and train the speech encoder and the adapter with contrastive loss,

$$\mathcal{L}_{\text{CTR}} = -\mathbb{E}_i \left[\log \frac{\exp(\text{sim}(W_i^s, W_i^t)/\tau)}{\sum_j \exp(\text{sim}(W_i^s, W_j^t)/\tau)} \right] \quad (7)$$

where τ is the temperature and $\text{sim}()$ is the cosine similarity function. LLM parameters are frozen during stage 1.

Stage 2. Finetuning for simultaneous translation. We finetune the entire model for simultaneous speech translation using wait- k -stride- n policy. Speech input is encoded into speech embeddings E^s . Then we concatenate E^s with embeddings of reference translation and feed them to LLM. Position indices of both speech embeddings and translation embeddings start with the same index and ascend separately, so that text generation during inference does not affect the positional embeddings of speech embeddings.

Then we randomly select $k \in K$ and mask out attentions from translation words with indices

¹<https://github.com/MontrealCorpusTools/Montreal-Forced-Aligner>

from in to $(i+1)n-1$ to speech segments $S_{>i+k}$ for each i , since these words are generated before speech segments $S_{>i+k}$ arrive during inference. Finally, we apply the cross-entropy loss to train the entire model,

$$\mathcal{L}_{\text{CE}} = -\mathbb{E}_i [\text{LLM}(T_i^* | E^s, T_{<i}^*, \text{mask})], \quad (8)$$

where T^* is the reference translation.

3.4 Efficient Simultaneous Inference

Figure 2 (c) illustrates how we conduct efficient simultaneous inference. FASST waits for k segments at the beginning and then start generating. Suppose now we have received S_1, S_2, \dots, S_i where $i \geq k$.

Incremental Speech Encoding The blockwise-causal mask of speech encoder allows us to use KV cache of previous speech segments to avoid recomputation. Let $H^s = (h_1^s, \dots, h_{l_i}^s)$ be input vectors of the attention. We group them into blocks $B_j = (h_{(j-1)b+1}^s, \dots, h_{jb}^s)$ where $1 \leq j \leq i$ and $i \cdot b = l_i$. The query, key and value matrices can be written as follows

$$Q = H^s M_Q = (B_1 M_Q, \dots, B_i M_Q) \quad (9)$$

$$K = H^s M_K = (B_1 M_K, \dots, B_i M_K) \quad (10)$$

$$V = H^s M_V = (B_1 M_V, \dots, B_i M_V) \quad (11)$$

Here the keys and values of previous segments $(B_1 M_K, \dots, B_{i-1} M_K)$ and $(B_1 M_V, \dots, B_{i-1} M_V)$ are stored in the KV cache. Now we only need the KV cache and the query $B_i M_Q$, key $B_i M_K$ and value $B_i M_V$ of the latest segment to compute its attention output,

$$O_i^s = \text{Softmax} \left(\frac{B_i M_Q K^T}{\sqrt{d}} \right) V. \quad (12)$$

This results in same output as running attention with full query, key and value matrices and a blockwise-causal mask. In this way, we reduce the time complexity of attention from $O(l_i d^2 + l_i^2 d)$ to $O(bd^2 + l_i bd)$. Here b is a constant while l_i increases with the longer speech input.

Adapting We store the speech encoder outputs of previous segments and concatenate them with encoder outputs of segment i . Then we pass them to the causal convolutional layers and the linear layer to obtain the speech embeddings $E_{\leq i}^s$.

290 **LLM Decoding with Consistency Mask** We par-
 291 tition speech embeddings $E_{\leq i}^s$ into E_1^s, \dots, E_i^s
 292 by speech segment. Following the wait- k -stride- n , the
 293 inputs to LLM are organized in the follow way

$$294 \quad I = E_1^s \oplus \dots \oplus E_k^s \oplus \text{Emb}(T_k) \oplus E_{k+1}^s \oplus \\
 295 \quad \text{Emb}(T_{k+1}) \oplus \dots \oplus E_i^s, \quad (13)$$

296 where $T_{1:k-1}$ are empty strings and T_j consists of
 297 n words for each $k \leq j < i$. Now we need to
 298 reuse KV cache of previous $i - 1$ speech segments
 299 and partial translations to compute LLM hidden
 300 states of i_{th} segment. Since speech embeddings are
 301 always ahead of text embeddings during training,
 302 we design a consistency mask to ensure speech seg-
 303 ments can only attend to speech segments before
 304 them.

305 Let $\delta(z)$ be indicator function that equals to 1
 306 if z_{th} position of input I belongs to text and 0
 307 otherwise. Define consistency mask M^c as follows,

$$308 \quad M_{z_Q, z_K}^c = \begin{cases} 0 & z_Q \geq z_K \text{ and } \delta(z_Q) \geq \delta(z_K) \\ -\infty & \text{otherwise} \end{cases} \quad (14)$$

309 Let $Q_i, K_i, V_i \in \mathbb{R}^{t_i \times d}$ be query, key and value
 310 matrices of segment i and $K_{< i}, V_{< i}$ be cached key
 311 and value matrices. We first concatenate K_i and
 312 V_i with cache to obtain $K_{\leq i}$ and $V_{\leq i}$. The atten-
 313 tion output of segment i can then be computed as
 314 follows

$$315 \quad O_i^t = \text{Softmax} \left(\frac{Q_i K_{\leq i}^T}{\sqrt{d}} + M_{-t_i, :}^c \right) V_{\leq i}. \quad (15)$$

316 After computing hidden states for speech seg-
 317 ment S_i , the LLM decodes n words autoregres-
 318 sively following the policy.

319 4 Experiment

320 4.1 Dataset

321 We conduct experiments on two language direc-
 322 tions of MuST-C v1.0 dataset (Di Gangi
 323 et al., 2019): English→Spanish (En-Es) and
 324 English→German (En-De). Each language direc-
 325 tion contains around 400 hours of audio recordings.
 326 The average duration of utterances is less than 10
 327 seconds. To simulate long speech scenarios, we
 328 concatenate adjacent utterances in the same TED
 329 talk so that each resulting utterance is around 30

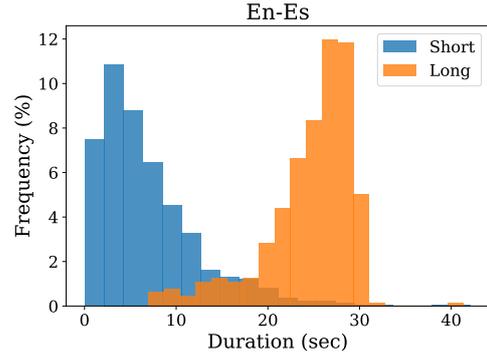


Figure 4: Duration distribution of MuST-C-Short and MuST-C-Long. The average duration of MuST-C-Short is around 5 seconds while that of MuST-C-Long is around 25 seconds.

seconds. We call the induced dataset as MuST-C-Long² and the original one as MuST-C-Short. The duration distribution of both datasets are shown in Figure 4.

334 4.2 Model Configurations

335 **Architecture** We initialize our speech encoder
 336 with wav2vec 2.0 large model³ (Baeviski et al.,
 337 2020) and our LLM with Llama2 7b Base⁴ (Tou-
 338 vron et al., 2023a). Wav2vec 2.0 large consists
 339 of a 7-layer convolutional feature extractor and a
 340 24-layer Transformer encoder with 1024 hidden
 341 units. The block size of speech encoder is set to
 342 50, i.e., around 1 second each block. The adapter
 343 connecting wav2vec 2.0 and Llama2 consists of
 344 two 1-D convolutional layers with kernel size 3,
 345 stride 2 and hidden size 1024 and a linear layer to
 346 project hidden size from 1024 to 4096 to match
 347 that of LLM embedding. Llama2 7b Base adopts
 348 a 32-layer Transformer decoder with hidden size
 349 4096. It uses a vocabulary of size 32000 and rotary
 350 positional embedding (Su et al., 2023).

351 **Training** We train our model with mixed MuST-
 352 C-Short and MuST-C-Long data. The input speech
 353 is raw 16-bit 16kHz mono-channel waveform. We
 354 filter out speech that is shorter than 320ms during
 355 training. The batch size of stage 1 is 16.7 minutes
 356 and that of stage 2 is 14 minutes. We use AdamW
 357 optimizer with cosine learning rate decay. The
 358 warmup steps of stage 1 is 25k and that of stage 2
 359 is 500 steps. The maximum learning rate of stage

²The manifest of MuST-C-Long will be released together with the code.

³https://dl.fbaipublicfiles.com/fairseq/wav2vec/wav2vec_vox_960h_pl.pt

⁴<https://huggingface.co/meta-llama/Llama-2-7b>

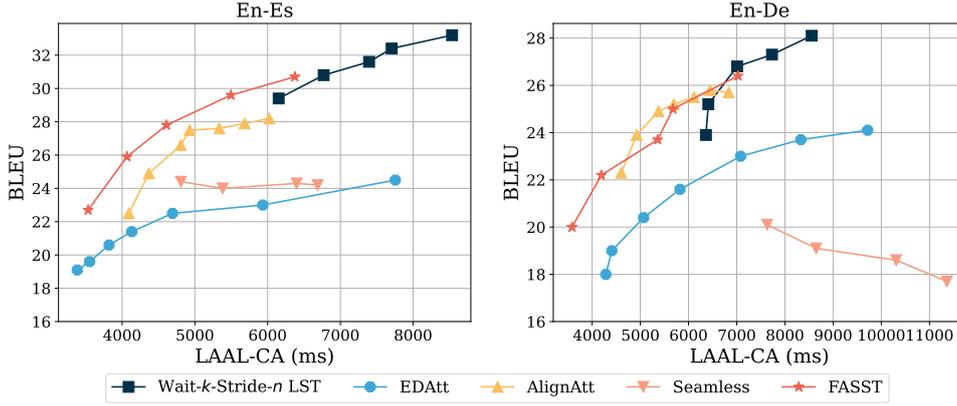


Figure 5: Quality-latency trade-off of FASST and baselines on English-Spanish and English-German direction. Quality is reflected by BLEU and latency is reflected by computation-aware length-adaptive average lagging (LAAL-CA). Given long speech input and large batch size, our model achieves overall the best quality-latency trade-off.

1 is $1e-4$ and that of stage 2 is $2e-5$. Gradients are clipped by the norm of 10. We train stage 1 for 500k steps and stage 2 for 1 epoch. We choose the checkpoint with the lowest dev loss. All our models are trained on 4 Nvidia A6000 GPUs with fp16. The temperature of WACO loss is 0.2 and we set $K = \{1, 2, 3, 4, 5, 100\}$, $n = 3$ for stage 2 training.

Inference We set speech segment size to 1 second to match the block size. We wait for $1 \leq k \leq 5$ segments at first. Then as each segment arrives, the speech encoder encodes its speech embedding incrementally and pass it to LLM, where LLM computes hidden states without recomputation and generates $n = 3$ words with greedy decoding as the partial translation.

4.3 Evaluation

We use SimulEval (Ma et al., 2020a) to evaluate our models and baselines. All models are evaluated on MuST-C-Long tst-COMMON with batch size of 8 during inference to simulate heavy workload. Since SimulEval does not support batching multiple instances, we duplicate each instance by 8 during model forwarding. We report SacreBLEU (Post, 2018) for translation quality and computation-aware length-adaptive average lagging (LAAL-CA) (Papi et al., 2022) for latency. All models are evaluated using a single A6000 GPU.

4.4 Baselines

Wait- k -Stride- n LST waits k fixed-length speech segments and translates n words every time (Ma et al., 2020b; Zeng et al., 2021). We run wait-

k -stride- n policy on a strong offline LLM-based model LST (Zhang et al., 2023b) trained on the same mixed data as FASST. LST has the Encoder-Adapter-LLM architecture similar to FASST but employs bidirectional speech encoder and requires recomputation every time a new speech segment arrives. We set $k \in \{1, 2, 3, 4, 5\}$, $n = 3$ and segment length 1 second to match the setting of FASST.

EDAtt is an attention-based adaptive policy (Papi et al., 2023a). It leverages the encoder-decoder attention of an offline ST model to decide when to emit partial translations. The intuition is that if the attention is focused on early audio frames, the current translation can be emitted since sufficient information has been received. We use the model checkpoint and settings provided by the authors.

AlignAtt is the current state-of-the-art (SOTA) method that extends EDAtt by explicitly generating audio-translation alignments from encoder-decoder attention (Papi et al., 2023b). While EDAtt emits based on attention scores directly, AlignAtt decides based on whether a predicted token aligns with the latest audio frames, providing a more interpretable latency control. We also use the model checkpoint and optimal settings provided by the authors of AlignAtt.

Seamless is a multilingual streaming speech translation system with efficient monotonic multi-head attention mechanism (Ma et al., 2023) to generate low latency translation (Communication et al., 2023). It computes target to source alignment using cross attention and writes translation if the

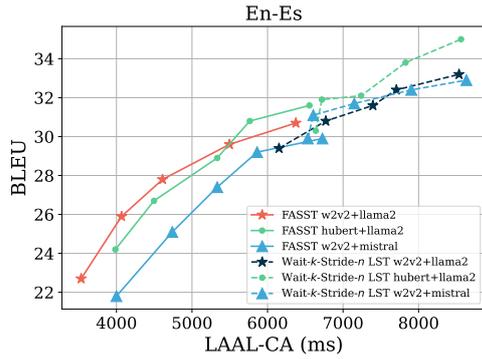


Figure 6: Ablation on the choice of pretrained speech encoder and LLM. We replace wav2vec 2.0 large with HuBERT large and Llama2 7B with Mistral 7B v0.3 base. FASST consistently has lower latency than Wait- k -Stride- n LST while maintaining an acceptable translation quality.

alignment probability is larger than a threshold. We vary the threshold in $[0.2, 0.4, 0.6, 0.8]$ to evaluate its quality-latency trade-off.

4.5 Main Results

Main results are shown in Figure 5. Our model achieves the best quality-latency trade-off for En-Es direction. Although wait- k -stride- n LST has a 2 BLEU score advantage at the latency of 8 seconds, its bidirectional encoding and inefficient use of KV cache prohibit it reaching latency smaller than 6 seconds. Comparing to EDAtt and AlignAtt which do not use LLM and has much less parameters (115M) than our model (7B), our model has similar computation aware latency while achieving a 1.5 BLEU score improvement. For En-De direction, FASST achieves competitive results to AlignAtt, with slightly better quality when latency is smaller than 4 seconds or larger than 6 seconds.

4.6 Ablation Studies

We conduct ablation studies to examine the impact of each component in our model.

Speech Encoder and LLM We replace wav2vec 2.0 large with HuBERT large (Hsu et al., 2021) and Llama2 7B base with Mistral 7B v0.3 base (Jiang et al., 2023) to examine whether FASST is sensitive to the choice of pretrained speech encoder and LLM. We also train Wait- k -Stride- n LST baseline with these configurations as a comparison. Results are shown in Figure 6. For all configurations, FASST has lower latency than the baseline. FASST with HuBERT results in the best quality

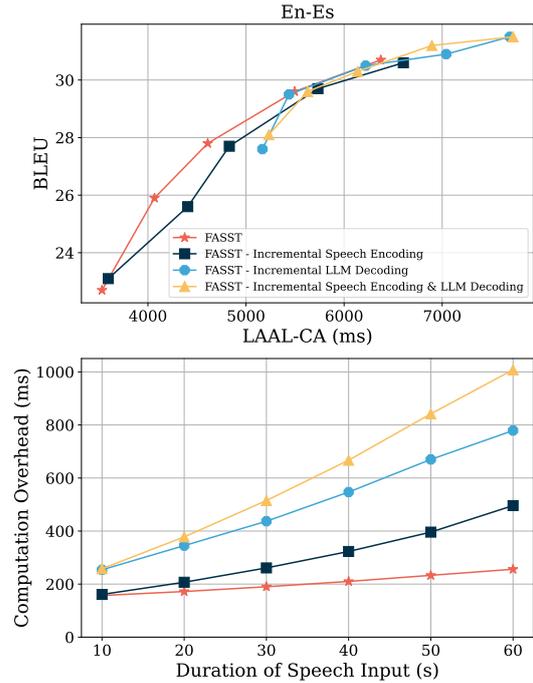


Figure 7: Ablation on incremental speech encoding and LLM decoding. Removing both incremental encoding and decoding can result in 4x larger latency for 60 seconds speech input.

when the latency is around 5.5~6.5 seconds and FASST with wav2vec 2.0 becomes the best when the latency is smaller than 5.5 seconds.

Incremental Encoding and Decoding We ablate the incremental speech encoding and the incremental LLM decoding to examine their impact. For encoding, we use the same architecture but recompute the entire speech encoder at each step. For decoding, we recompute the entire LLM hidden states given updated speech input and then incrementally decode the translation as each speech segment arrives. This also provides translation tokens with more context since they can attend to speech embeddings appear after them. Results are shown in Figure 7. Incremental encoding of speech encoder reduces the computational latency consistently by at least 200ms compared to recomputing encoder. Recomputing LLM does improve translation quality (1 ~ 5 BLEU), but also introduces significant computation overhead (~ 1.5 second), making it inferior for real-time application.

We also plot the computation cost of each read/write step for each variant in Figure 7 with wait-2-stride-3 policy. FASST scales the best with the speech input length and reduces the overhead by at most 4x compared to the one without incremen-

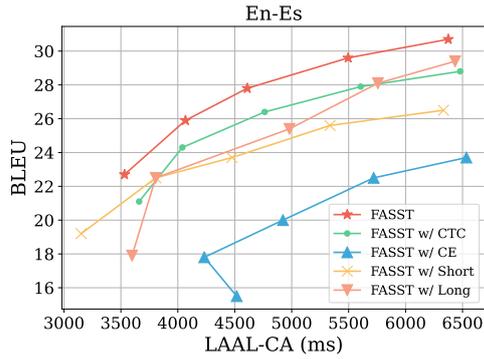


Figure 8: Ablation on the training strategy of FASST. We train stage 1 WACO loss with CTC and cross entropy (CE) loss and also change the training data to only MuST-C-Short and only MuST-C-Long. WACO and mixed-data training achieves the best performance.

tal encoding and decoding. Removing incremental decoding results in larger additional computation overhead compared to removing incremental encoding since LLM has far more parameters than the speech encoder.

Speech-Text Alignment Training We replace WACO loss with CTC loss (Graves et al., 2006) and cross entropy loss in stage 1 to examine its impact on model performance. CTC loss aligns speech and text embeddings by optimizing all possible alignment paths. For cross entropy loss, we pass the speech embeddings to LLM and optimize the cross entropy loss with LLM parameters frozen. As shown in Figure 8, WACO loss consistently outperforms CTC and cross entropy by at least 1 BLEU score at the same latency.

Mixing MuST-C-Long and Short We train our model separately with only MuST-C-Short and only MuST-C-Long for the same number of epochs. As shown in Figure 8, the model trained with both long and short data outperforms the one trained with short data by up to 4 BLEU points. Though we are using LLM, the length extrapolation is still unsatisfactory. Training with long data improves the quality compared to short data at high latency, but still outperformed by mixed-data training.

4.7 Generalizability to Other Policy

We have demonstrated that our method works with wait- k -stride- n policy. However, plenty of policies other than wait- k and its variants have been developed to conduct simultaneous translation. Here we apply our method to hold- n policy (Liu et al., 2020) to exemplify how our method works on a different

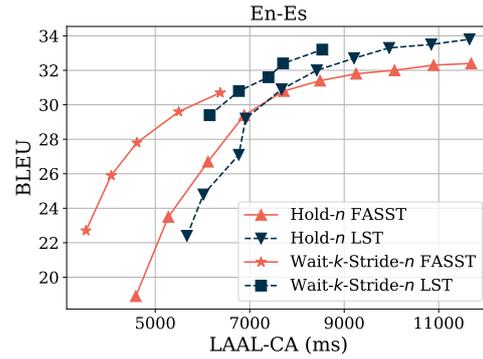


Figure 9: Quality-latency trade-off of FASST when applied to hold- n policy. We observe less improvement with respect to LST than wait- k -stride- n policy.

policy and in the meanwhile explain the factors that influence the effectiveness of our method.

Hold- n policy selects the best hypothesis after each speech segment arrives, discard the last n tokens from it and outputs the rest as the partial translation. Since hold- n is a test-time policy for offline ST model, we train an offline version of our model by replacing stage 2 finetuning with standard offline finetuning using cross entropy loss but keep blockwise-causal encoding. During inference, we still conduct the same incremental encoding and decoding.

As shown in Figure 9, our method has less advantage when applied to hold- n comparing to wait- k -stride-3. The major difference between two policies in terms of computation is that hold- n policy spends more time on autoregressive decoding since it decodes more tokens each time. On average hold- n policy generates more than 6 words each step while wait- k -stride-3 generates at most 3 words. FASST accelerates the encoding of existing features, but for policies like hold- n that involve heavy autoregressive decoding the advantage of our method gets marginalized.

5 Conclusion

In this work, we introduce FASST, a fast LLM-based simultaneous speech translation model. FASST consists of blockwise-causal speech encoding, incremental LLM decoding with consistency mask, and a novel 2-stage training strategy. Experiments on MuST-C dataset show that FASST significantly reduce computation overhead while maintaining its translation quality. Our generalization study shows that policies that spend more time on encoding than decoding benefit more from FASST.

550 Limitations

- 551 • There might be data leakage since LLM is trained
552 on vast amount of text data, so we cannot guar-
553 antee LLM does not see the test translation data
554 during pretraining.
- 555 • FASST is only tested on two language directions
556 instead of all 8 language directions of MuST-C,
557 so its generalizability to other language direc-
558 tions is unknown.
- 559 • There is still a quality gap between blockwise-
560 causal speech encoding and bidirectional speech
561 encoding. It is unclear how to further close the
562 gap.
- 563 • We only explore one LLM-ST architecture in the
564 paper and we cannot guarantee that FASST or its
565 idea works on other architectures.

566 References

567 Milind Agarwal, Sweta Agrawal, Antonios Anasta-
568 sopoulos, Luisa Bentivogli, Ondřej Bojar, Claudia
569 Borg, Marine Carpuat, Roldano Cattoni, Mauro Cet-
570 tolo, Mingda Chen, William Chen, Khalid Choukri,
571 Alexandra Chronopoulou, Anna Currey, Thierry De-
572 clerck, Qianqian Dong, Kevin Duh, Yannick Es-
573 tève, Marcello Federico, Souhir Gahbiche, Barry
574 Haddow, Benjamin Hsu, Phu Mon Htut, Hirofumi
575 Inaguma, Dávid Javorský, John Judge, Yasumasa
576 Kano, Tom Ko, Rishu Kumar, Pengwei Li, Xutai Ma,
577 Prashant Mathur, Evgeny Matusov, Paul McNamee,
578 John P. McCrae, Kenton Murray, Maria Nadejde,
579 Satoshi Nakamura, Matteo Negri, Ha Nguyen, Jan
580 Niehues, Xing Niu, Atul Kr. Ojha, John E. Ortega,
581 Proyag Pal, Juan Pino, Lonneke van der Plas, Peter
582 Polák, Elijah Rippeth, Elizabeth Salesky, Jiatong Shi,
583 Matthias Sperber, Sebastian Stüker, Katsuhito Su-
584 doh, Yun Tang, Brian Thompson, Kevin Tran, Marco
585 Turchi, Alex Waibel, Mingxuan Wang, Shinji Watan-
586 abe, and Rodolfo Zevallos. 2023. **FINDINGS OF
587 THE IWSLT 2023 EVALUATION CAMPAIGN**. In
588 *Proceedings of the 20th International Conference on
589 Spoken Language Translation (IWSLT 2023)*, pages
590 1–61, Toronto, Canada (in-person and online). Asso-
591 ciation for Computational Linguistics.

592 Alexei Baevski, Henry Zhou, Abdelrahman Mohamed,
593 and Michael Auli. 2020. **wav2vec 2.0: A framework
594 for self-supervised learning of speech representations**.
595 *Preprint*, arXiv:2006.11477.

596 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie
597 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind
598 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
599 Askeel, Sandhini Agarwal, Ariel Herbert-Voss,
600 Gretchen Krueger, Tom Henighan, Rewon Child,
601 Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,
602 Clemens Winter, Christopher Hesse, Mark Chen,
603 Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin

604 Chess, Jack Clark, Christopher Berner, Sam Mc-
605 Candlish, Alec Radford, Ilya Sutskever, and Dario
606 Amodei. 2020. **Language models are few-shot learn-
607 ers**. *Preprint*, arXiv:2005.14165.

Junkun Chen, Mingbo Ma, Renjie Zheng, and Liang
608 Huang. 2021. **Direct simultaneous speech-to-text
609 translation assisted by synchronized streaming ASR**.
610 In *Findings of the Association for Computational
611 Linguistics: ACL-IJCNLP 2021*, pages 4618–4624,
612 Online. Association for Computational Linguistics.
613

Seamless Communication, Loïc Barrault, Yu-An Chung,
614 Mariano Coria Meglioli, David Dale, Ning Dong,
615 Mark Duppenthaler, Paul-Ambroise Duquenne,
616 Brian Ellis, Hady Elsahar, Justin Haaheim, John Hoff-
617 man, Min-Jae Hwang, Hirofumi Inaguma, Christo-
618 pher Klaiber, Ilia Kulikov, Pengwei Li, Daniel Licht,
619 Jean Maillard, Ruslan Mavlyutov, Alice Rakotoari-
620 son, Kaushik Ram Sadagopan, Abinеш Ramakr-
621 ishnan, Tuan Tran, Guillaume Wenzek, Yilin Yang,
622 Ethan Ye, Ivan Evtimov, Pierre Fernandez, Cynthia
623 Gao, Prangthip Hansanti, Elahe Kalbassi, Amanda
624 Kallet, Artyom Kozhevnikov, Gabriel Mejia Gonza-
625 lez, Robin San Roman, Christophe Touret, Corinne
626 Wong, Carleigh Wood, Bokai Yu, Pierre Andrews,
627 Can Balioglu, Peng-Jen Chen, Marta R. Costa-jussà,
628 Maha Elbayad, Hongyu Gong, Francisco Guzmán,
629 Kevin Heffernan, Somya Jain, Justine Kao, Ann
630 Lee, Xutai Ma, Alex Mourachko, Benjamin Pello-
631 quin, Juan Pino, Sravya Popuri, Christophe Ropers,
632 Safiyah Saleem, Holger Schwenk, Anna Sun, Paden
633 Tomasello, Changan Wang, Jeff Wang, Skyler Wang,
634 and Mary Williamson. 2023. **Seamless: Multilin-
635 gual expressive and streaming speech translation**.
636 *Preprint*, arXiv:2312.05187.
637

Tri Dao. 2024. **Flashattention-2: Faster attention with
638 better parallelism and work partitioning**. In *The
639 Twelfth International Conference on Learning Repre-
640 sentations*.
641

Mattia A. Di Gangi, Roldano Cattoni, Luisa Bentivogli,
642 Matteo Negri, and Marco Turchi. 2019. **MuST-C: a
643 Multilingual Speech Translation Corpus**. In *Proceed-
644 ings of the 2019 Conference of the North American
645 Chapter of the Association for Computational Lin-
646 guistics: Human Language Technologies, Volume 1
647 (Long and Short Papers)*, pages 2012–2017, Min-
648 neapolis, Minnesota. Association for Computational
649 Linguistics.
650

Qian Dong, Yaoming Zhu, Mingxuan Wang, and Lei
651 Li. 2022. **Learning when to translate for streaming
652 speech**. In *Proceedings of the 60th Annual Meet-
653 ings of the Association for Computational Linguistics
654 (Volume 1: Long Papers)*, pages 680–694, Dublin,
655 Ireland. Association for Computational Linguistics.
656

Alex Graves, Santiago Fernández, Faustino Gomez, and
657 Jürgen Schmidhuber. 2006. Connectionist temporal
658 classification: labelling unsegmented sequence data
659 with recurrent neural networks. In *Proceedings of the
660 23rd international conference on Machine learning*,
661 pages 369–376.
662

663	Shoutao Guo, Shaolei Zhang, Zhengrui Ma, Min Zhang, and Yang Feng. 2024. Sillm: Large language models for simultaneous machine translation . <i>Preprint</i> , arXiv:2402.13036.	718
664		719
665		720
666		721
667	Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed. 2021. Hubert: Self-supervised speech representation learning by masked prediction of hidden units . <i>IEEE/ACM Trans. Audio, Speech and Lang. Proc.</i> , 29:3451–3460.	722
668		723
669		
670		724
671		725
672		726
673	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 . <i>Preprint</i> , arXiv:2312.13585.	727
674		728
675		729
676		
677	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L��lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth��e Lacroix, and William El Sayed. 2023. Mistral 7b . <i>Preprint</i> , arXiv:2310.06825.	730
678		731
679		732
680		733
681		734
682		735
683		736
684		
685	Roman Koshkin, Katsuhito Sudoh, and Satoshi Nakamura. 2024. Transllama: Llm-based simultaneous translation system . <i>Preprint</i> , arXiv:2402.04636.	737
686		738
687		739
688	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In <i>Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles</i> .	740
689		741
690		742
691		743
692		744
693		745
694		746
695	Danni Liu, Gerasimos Spanakis, and Jan Niehues. 2020. Low-Latency Sequence-to-Sequence Speech Recognition and Translation by Partial Hypothesis Selection . In <i>Proc. Interspeech 2020</i> , pages 3620–3624.	747
696		
697		748
698		749
699	Xutai Ma, Mohammad Javad Dousti, Changhan Wang, Jiatao Gu, and Juan Pino. 2020a. SIMULEVAL: An evaluation toolkit for simultaneous translation . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 144–150, Online. Association for Computational Linguistics.	750
700		751
701		752
702		753
703		754
704		755
705		756
706	Xutai Ma, Juan Pino, and Philipp Koehn. 2020b. SimulMT to SimulST: Adapting simultaneous text translation to end-to-end simultaneous speech translation . In <i>Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing</i> , pages 582–587, Suzhou, China. Association for Computational Linguistics.	757
707		758
708		759
709		760
710		761
711		762
712		763
713		764
714		765
715	Xutai Ma, Anna Sun, Siqi Ouyang, Hirofumi Inaguma, and Paden Tomasello. 2023. Efficient monotonic multihead attention . <i>Preprint</i> , arXiv:2312.04515.	766
716		767
717		768
		769
		770
		771
		772
	Xutai Ma, Yongqiang Wang, Mohammad Javad Dousti, Philipp Koehn, and Juan Pino. 2020c. Streaming simultaneous speech translation with augmented memory transformer . <i>Preprint</i> , arXiv:2011.00033.	
	R OpenAI. 2023. Gpt-4 technical report . arxiv 2303.08774. <i>View in Article</i> , 2:13.	
	Siqi Ouyang, Rong Ye, and Lei Li. 2023. WACO: Word-aligned contrastive learning for speech translation . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3891–3907, Toronto, Canada. Association for Computational Linguistics.	
	Sara Papi, Marco Gaido, Matteo Negri, and Marco Turchi. 2022. Over-generation cannot be rewarded: Length-adaptive average lagging for simultaneous speech translation . In <i>Proceedings of the Third Workshop on Automatic Simultaneous Translation</i> , pages 12–17, Online. Association for Computational Linguistics.	
	Sara Papi, Matteo Negri, and Marco Turchi. 2023a. Attention as a guide for simultaneous speech translation . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 13340–13356, Toronto, Canada. Association for Computational Linguistics.	
	Sara Papi, Marco Turchi, and Matteo Negri. 2023b. Alignatt: Using attention-based audio-translation alignments as a guide for simultaneous speech translation . In <i>INTERSPEECH 2023</i> , interspeech_2023. ISCA.	
	Peter Pol��k, Brian Yan, Shinji Watanabe, Alex Waibel, and Ondr��j Bojar. 2023. Incremental Blockwise Beam Search for Simultaneous Speech Translation with Controllable Quality-Latency Tradeoff . In <i>Proc. INTERSPEECH 2023</i> , pages 3979–3983.	
	Reiner Pope, Sholto Douglas, Aakanksha Chowdhery, Jacob Devlin, James Bradbury, Anselm Levskaya, Jonathan Heek, Kefan Xiao, Shivani Agrawal, and Jeff Dean. 2022. Efficiently scaling transformer inference . <i>Preprint</i> , arXiv:2211.05102.	
	Matt Post. 2018. A call for clarity in reporting BLEU scores . In <i>Proceedings of the Third Conference on Machine Translation: Research Papers</i> , pages 186–191, Brussels, Belgium. Association for Computational Linguistics.	
	Matthew Raffel, Victor Agostinelli, and Lihong Chen. 2024. Simultaneous masking, not prompting optimization: A paradigm shift in fine-tuning llms for simultaneous translation . <i>Preprint</i> , arXiv:2405.10443.	
	Matthew Raffel and Lihong Chen. 2023. Implicit memory transformer for computationally efficient simultaneous speech translation . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 12900–12907, Toronto, Canada. Association for Computational Linguistics.	

773	Matthew Raffel, Drew Penney, and Lizhong Chen. 2023.	lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,	833
774	Shiftable context: addressing training-inference con-	Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,	834
775	text mismatch in simultaneous speech translation.	Melanie Kambadur, Sharan Narang, Aurelien Ro-	835
776	In <i>International Conference on Machine Learning</i> ,	driguez, Robert Stojnic, Sergey Edunov, and Thomas	836
777	pages 28519–28530. PMLR.	Scialom. 2023b. Llama 2: Open foundation and	837
		fine-tuned chat models . <i>Preprint</i> , arXiv:2307.09288.	838
778	Yi Ren, Jinglin Liu, Xu Tan, Chen Zhang, Tao Qin,	David Vilar, Markus Freitag, Colin Cherry, Jiaming Luo,	839
779	Zhou Zhao, and Tie-Yan Liu. 2020. SimulSpeech:	Viresh Ratnakar, and George Foster. 2023. Prompt-	840
780	End-to-end simultaneous speech to text translation .	ing PaLM for translation: Assessing strategies and	841
781	In <i>Proceedings of the 58th Annual Meeting of the As-</i>	performance . In <i>Proceedings of the 61st Annual</i>	842
782	<i>sociation for Computational Linguistics</i> , pages 3787–	<i>Meeting of the Association for Computational Lin-</i>	843
783	3796, Online. Association for Computational Lin-	<i>guistics (Volume 1: Long Papers)</i> , pages 15406–	844
784	guistics.	15427, Toronto, Canada. Association for Computa-	845
785	Sascha Rothe, Shashi Narayan, and Aliaksei Severyn.	tional Linguistics.	846
786	2020. Leveraging pre-trained checkpoints for se-		
787	quence generation tasks . <i>Transactions of the Associ-</i>	Chunyang Wu, Yongqiang Wang, Yangyang Shi, Ching-	847
788	<i>ation for Computational Linguistics</i> , 8:264–280.	Feng Yeh, and Frank Zhang. 2020. Stream-	848
789	Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha,	ing transformer-based acoustic models using self-	849
790	Bo Wen, and Yunfeng Liu. 2023. Roformer: En-	attention with augmented memory . <i>Preprint</i> ,	850
791	hanced transformer with rotary position embedding .	arXiv:2005.08042.	851
792	<i>Preprint</i> , arXiv:2104.09864.		
793	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	Xueqing Wu, Lewen Wang, Yingce Xia, Weiqing Liu,	852
794	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	Lijun Wu, Shufang Xie, Tao Qin, and Tie-Yan Liu.	853
795	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	2021. Temporally correlated task scheduling for se-	854
796	Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton	quence learning. In <i>International Conference on</i>	855
797	Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,	<i>Machine Learning</i> , pages 11274–11284. PMLR.	856
798	Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller,		
799	Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-	Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Has-	857
800	thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan	san Awadalla. 2023. A paradigm shift in machine	858
801	Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,	translation: Boosting translation performance of	859
802	Isabel Kloumann, Artem Korenev, Punit Singh Koura,	large language models . <i>Preprint</i> , arXiv:2309.11674.	860
803	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-		
804	ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-	Wen Yang, Chong Li, Jiajun Zhang, and Chengqing	861
805	tinet, Todor Mihaylov, Pushkar Mishra, Igor Moly-	Zong. 2023. Bigtranslate: Augmenting large lan-	862
806	bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-	guage models with multilingual translation capability	863
807	stein, Rashi Rungta, Kalyan Saladi, Alan Schelten,	over 100 languages . <i>Preprint</i> , arXiv:2305.18098.	864
808	Ruan Silva, Eric Michael Smith, Ranjan Subrama-		
809	nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-	Xingshan Zeng, Liangyou Li, and Qun Liu. 2021. Real-	865
810	lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,	TranS: End-to-end simultaneous speech translation	866
811	Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,	with convolutional weighted-shrinking transformer .	867
812	Melanie Kambadur, Sharan Narang, Aurelien Ro-	In <i>Findings of the Association for Computational</i>	868
813	driguez, Robert Stojnic, Sergey Edunov, and Thomas	<i>Linguistics: ACL-IJCNLP 2021</i> , pages 2461–2474,	869
814	Scialom. 2023a. Llama 2: Open foundation and fine-	Online. Association for Computational Linguistics.	870
815	tuned chat models . <i>Preprint</i> , arXiv:2307.09288.		
816	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	Xingshan Zeng, Liangyou Li, and Qun Liu. 2023. Ada-	871
817	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	TranS: Adapting with boundary-based shrinking for	872
818	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	end-to-end speech translation . In <i>Findings of the</i>	873
819	Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton	<i>Association for Computational Linguistics: EMNLP</i>	874
820	Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,	2023, pages 2353–2361, Singapore. Association for	875
821	Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller,	Computational Linguistics.	876
822	Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-		
823	thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan	Biao Zhang, Barry Haddow, and Alexandra Birch.	877
824	Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,	2023a. Prompting large language model for ma-	878
825	Isabel Kloumann, Artem Korenev, Punit Singh Koura,	chine translation: a case study. In <i>Proceedings of</i>	879
826	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-	<i>the 40th International Conference on Machine Learn-</i>	880
827	ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-	<i>ing</i> , ICML’23. JMLR.org.	881
828	tinet, Todor Mihaylov, Pushkar Mishra, Igor Moly-		
829	bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-	Hao Zhang, Nianwen Si, Yaqi Chen, Wenlin Zhang,	882
830	stein, Rashi Rungta, Kalyan Saladi, Alan Schelten,	Xukui Yang, Dan Qu, and Xiaolin Jiao. 2023b. Tun-	883
831	Ruan Silva, Eric Michael Smith, Ranjan Subrama-	ing large language model for end-to-end speech trans-	884
832	nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-	lation . <i>Preprint</i> , arXiv:2310.02050.	885

886 Linlin Zhang, Kai Fan, Jiajun Bu, and Zhongqiang
887 Huang. 2023c. [Training simultaneous speech trans-](#)
888 [lation with robust and random wait-k-tokens strat-](#)
889 [egy](#). In *Proceedings of the 2023 Conference on Em-*
890 *pirical Methods in Natural Language Processing*,
891 pages 7814–7831, Singapore. Association for Com-
892 putational Linguistics.

893 Shaolei Zhang, Qingkai Fang, Zhuocheng Zhang, Zhen-
894 grui Ma, Yan Zhou, Langlin Huang, Mengyu Bu,
895 Shangtong Gui, Yunji Chen, Xilin Chen, and Yang
896 Feng. 2023d. [Bayling: Bridging cross-lingual align-](#)
897 [ment and instruction following through interactive](#)
898 [translation for large language models](#). *Preprint*,
899 arXiv:2306.10968.

900 Shaolei Zhang and Yang Feng. 2022. [Information-](#)
901 [transport-based policy for simultaneous translation](#).
902 In *Proceedings of the 2022 Conference on Empirical*
903 *Methods in Natural Language Processing*, pages 992–
904 1013, Abu Dhabi, United Arab Emirates. Association
905 for Computational Linguistics.

906 Shaolei Zhang and Yang Feng. 2023. [End-to-end simul-](#)
907 [taneous speech translation with differentiable seg-](#)
908 [mentation](#). In *Findings of the Association for Com-*
909 *putational Linguistics: ACL 2023*, pages 7659–7680,
910 Toronto, Canada. Association for Computational Lin-
911 guistics.