

Overcoming Latency Bottlenecks in On-Device Speech Translation: A Cascaded Approach with Alignment-Based Streaming MT

Zeeshan Ahmed, Frank Seide, Niko Moritz, Ju Lin, Ruiming Xie, Simone Merello, Zhe Liu and Christian Fuegen
Meta AI, USA

Abstract—This paper tackles several challenges that arise when integrating Automatic Speech Recognition (ASR) and Machine Translation (MT) for real-time, on-device streaming speech translation. Although state-of-the-art ASR systems based on Recurrent Neural Network Transducers (RNN-T) can perform real-time transcription, achieving streaming translation in real-time remains a significant challenge. To address this issue, we propose a simultaneous translation approach that effectively balances translation quality and latency. We also investigate efficient integration of ASR and MT, leveraging linguistic cues generated by the ASR system to manage context and utilizing efficient beam-search pruning techniques such as time-out and forced finalization to maintain system’s real-time factor. We apply our approach to an on-device bilingual conversational speech translation and demonstrate that our techniques outperform baselines in terms of latency and quality. Notably, our technique narrows the quality gap with non-streaming translation systems, paving the way for more accurate and efficient real-time speech translation.

Index Terms—streaming speech translation, ASR, MT, Simultaneous Machine Translation

I. INTRODUCTION

The traditional approach to speech translation involves a pipeline that combines an Automatic Speech Recognition (ASR) system with a Machine Translation (MT) system. However, current state-of-the-art cascade-based Speech Translation (ST) systems are for offline use, where the entire audio input is available in advance. In contrast, streaming setups require real-time processing of unbounded audio inputs, which poses significant challenges for cascade-based ST systems [1], [2]. Conventional MT systems are the bottleneck for the real-time streaming translation set-up, as they are typically trained on sentence-aligned pairs of text. This creates a mismatch between the training data, which consists of short sentences with only a few hundred tokens, and the inference conditions (in a streaming environment), where thousands of tokens may be processed in live sessions. To address this issue, segmentation models are often used to divide the incoming text stream into manageable chunks, or segments, that can be translated independently by the MT system. However, translating the segment in isolation poses a significant risk to translation quality [3]–[5].

Simultaneous Machine Translation (SMT) [6]–[10], on the other hand, is better suited for streaming speech translation scenarios. The SMT initiates the translation process while still

receiving the input. Conventional MT relies on an encoder-decoder architecture. The decoder uses an attention module that attends to all encoder embeddings. For SMT, both the encoder and attention module must be adaptive, allowing the encoder to operate incrementally and the cross-attention in the decoder to handle variable-length encoder output sequences. In this setup, the SMT model is fed one token at a time, and at each decoding step, it decides whether it wants to wait for another token or generate a translation. The SMT systems face a trade-off between quality and latency. Prioritizing higher quality increases latency (with conventional non-streaming MT systems still representing the upper bound in quality), while reducing latency can impact translation quality. Therefore, SMT systems must carefully balance this trade-off for practical applications. Another challenge for SMT is that, despite its ability to translate input in real-time, it still requires frequent context flushing to effectively manage context and resources, particularly in on-device scenarios.

In this paper, we propose a couple of advances in cascade-based speech translation. 1) We present a novel approach to SMT called ALiGnment BAsed STReaming Machine Translation (ALiBaStr-MT) which efficiently balances the quality/latency trade-off in on-device streaming speech translation. 2) We develop a robust multi-channel multi-talker ASR model that can recognize multiple languages, distinguish between different speakers and suppress side-talk/background noise. 3) Despite operating in a streaming fashion, the SMT architecture needs to keep all the previous encoder/decoder states saved in memory (kv-cached) for translation of future input. When used in an on-device setting, the memory is limited and the context needs frequent flushing. We propose incorporating segmentation cues in ASR to detect the segment boundaries and flush the states based on detected segments.

II. BACKGROUND

Cascade-based and end-to-end are two major paradigms for speech translation. Cascade-based systems use separate ASR and MT models [1], [2], while end-to-end systems translate directly from source speech to target text using a single model [11]. End-to-end systems offer advantages in terms of latency and model size, while cascade-based systems still demonstrate superior translation accuracy in many cases with well trained ASR and MT models. Furthermore, the decoupling

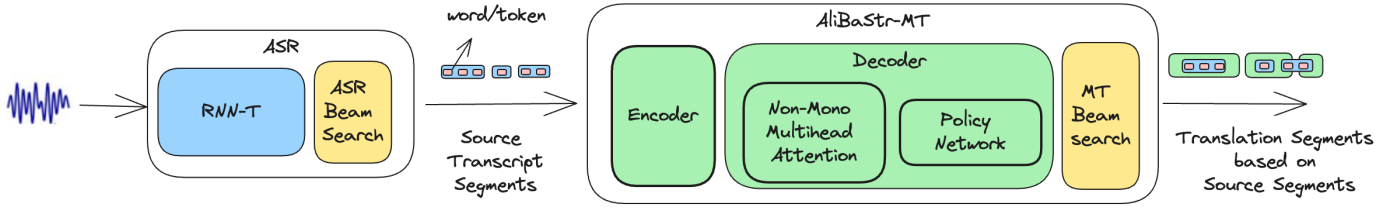


Fig. 1. On-device Streaming Speech Translation Architecture using Cascade of ASR and MT. The AliBaStr-MT system processes input segments of varying sizes and dynamically decides whether to generate a translation or await additional input, guided by its read/write policy module.

of ASR and MT modules also allows independent training and optimization that favors low resource language pairs [12] for which bilingual end-to-end audio training data are scarce. Incorporating streaming capabilities into these systems is a challenging task that has received increasing attention in recent years.

Although end-to-end models [13]–[16] are more favorable for streaming speech translation, cascade-based models [17] face considerable challenges. For instance, accurately identifying the correct segment of the ASR transcript for translation is a challenge. ASR systems often segment audio based on acoustic features such as pauses, which may not correspond to natural sentence boundaries [18]. This mismatch can lead to segments that are too long (spanning multiple sentences) or too short (containing only fragments), negatively affecting the accuracy of subsequent MT. To overcome this problem, researchers have explored re-segmenting ASR transcripts using punctuation restoration and language models to improve sentence boundaries and enhance MT performance [18], [19]. In this work, we incorporate the prediction of punctuation directly into the RNN-T model and employ punctuation for segment boundary detection (mostly sentence) when context flushing is needed.

Another challenge arises when attempting to translate these segments in real-time. Conventional MT systems are designed to translate complete sentences, and are not well-suited for translating smaller segments, making them impractical for real-time streaming translation. Simultaneous MT is particularly beneficial for the real-time speech translation scenario.

SMT systems can be broadly classified into fixed policy and adaptive policy systems. In the fixed policy system, the model uses the predefined heuristics and rules to generate the translation. These heuristics remain constant throughout the translation. Fixed policy systems such as wait- k or prefix-to-prefix translation systems by [20], [21] are simple and yet effective methods for streaming translation. However, the value of k in wait- k is not a learnable parameter of the model and is set beforehand, making it difficult to adjust based on the input and so far generated translation. The k is also language dependent for different translation directions, making it difficult to build a multilingual streaming translation model with wait- k . [22] on the other hand, tries to fix the quality gap in simultaneous translation models by introducing a look-ahead feature and the curriculum learning during training.

The adaptive methods, on the other hand, dynamically adjust

the read/write decision based on the input and previously generated translation. The adaptive method includes the ones where non-streaming MT model is kept fixed/frozen and a heuristics or learnable module is added to make read/write decisions e.g. the confidence threshold (What-if-*) method of [23] and reinforcement learning approaches by [24], [25] use a pretrained non-streaming model which is not modified during the training. [26] also proposed an adaptive system without attention. The model keeps the single context vector to reduce computational and memory footprint. In addition to the target tokens, the model can also generate the empty symbol to delay the translation. However, the model requires external supervision in the form of alignment between the source and target words to train the system. [27] compares the full sentence translation with the partial translation and builds up a ground-truth action sequence to train the read/write module. [28] explore the application of attention scores to create adaptive policies that guide the inference process.

Most of the recent work on adaptive policy methods is based on monotonic attention. The monotonic attention in essence learns a dynamic read/write policy which varies over the course of translation. There are various flavors available for monotonic attention, e.g. Monotonic Hard Attention [10] where the model only attends to the latest encoder state causing improvements in the latency and runtime while making a compromise on the translation quality. Monotonic Chunkwise Attention (MoChA) [9] on the other hand, attends to the segment of previous encoder states to overcome the quality gap. Monotonic Infinite Lookback (MILK) [8] further improves quality by attending to the entire previous encoder states. These flavors can also be combined with their multi-head version called Monotonic Multi-head Attention (MMA) [6], [7]. One of the drawbacks of the MA is that the read/write policy module is tightly embedded within the decoder layers. If there are n decoder layers, there will be n policy modules. Whereas, only the last policy module is used for the decision.

III. METHOD

Figure 1 depicts the overall architecture of our system. The system is composed of an streaming RNN-T based multi-channel multi-lingual ASR model with a beam-search module to generate the ASR transcripts. Translations for ASR transcripts are generated using the AliBaStr-MT system based on encoder/decoder architecture with an SMT specific beam-search component.

A. Automatic Speech Recognition (ASR)

Our ASR system is similar to [29], [30] which is based on RNN-T. This architecture is composed of three main components: an encoder, a prediction network, and a joiner network. The primary objective of the transducer model is to generate a sequence of labels $Y = (y_1, \dots, y_L)$ of length L . These labels can represent words or word pieces, derived from an input sequence $X = (x_1, \dots, x_S)$ of length S , which typically consists of acoustic features such as Mel-spectral energies.

Our ASR model is designed to effectively capture and transcribe conversations in various settings. It utilizes a multi-channel approach, leveraging microphone arrays placed at different locations on the device to robustly identify the change in speaker, distinguish between conversation partners' speech and suppress side-talk. Additionally, the model is streaming multi-lingual, automatically detecting language changes within an audio stream, whether from the same or multiple speakers, and generating corresponding transcripts. This functionality is achieved through serialized output training (SOT) and a multichannel front-end, as described in [29]. During training, the RNN-T model is also trained to emit four punctuation marks (‘,’ , ‘?’ , ‘!’ , ‘.’) to identify segment boundaries. Finally, we employ beam search decoding on top of RNN-T to generate accurate ASR transcripts, which are then fed into the SMT system for translation.

B. Simultaneous MT

We enable streaming translation ability in a pretrained non-streaming encoder/decoder model using an additional read/write policy module. Given that the quality of the non-streaming model is currently the upper bound for the streaming model, we propose to learn the read/write policy using the alignment computed from the non-monotonic attention in a supervised manner. We get these supervised labels called pseudo-labels automatically from attention weights. The overall architecture of our ALIGNment BASEd STreaming Machine Translation (AliBaStr-MT) is shown in Figure 2.

1) *Non-Streaming Full Sentence Translation Model:* The non-streaming model is trained on full sentence context using the following formulation [31]. Formally, a model based on standard encoder-decoder architecture transforms an input $x = \{x_1, \dots, x_{|x|}\}$ into an output sequence $y = \{y_1, \dots, y_{|y|}\}$ using the following

$$h_j = \text{Encoder}(x_i) \quad (1)$$

$$s_i = \text{Decoder}(y_{i-1}, s_{i-1}, c_i) \quad (2)$$

$$y_i = \text{Output}(s_i, c_i) \quad (3)$$

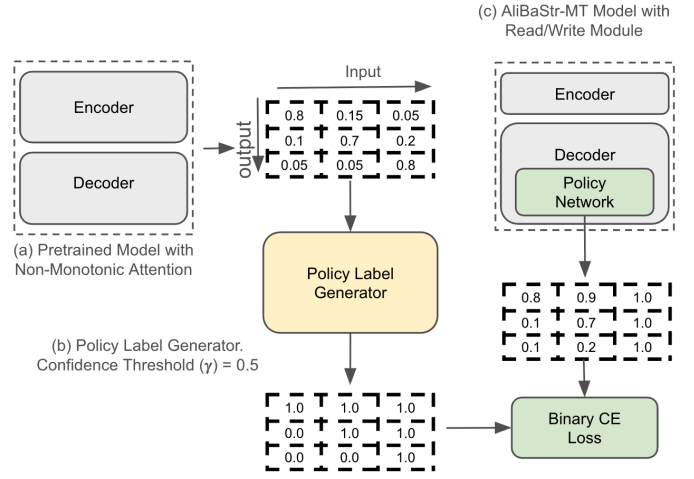


Fig. 2. Read/Write Policy learning in AliBaStr-MT using a pretrained non-streaming encoder-decoder model.

The context vector c_i is computed using a soft attention as follows.

$$e_{i,j} = \text{Energy}(s_{i-1}, h_j) \quad (4)$$

$$\alpha_{i,j} = \frac{\exp(e_{i,j})}{\sum_{k=1}^T \exp(e_{i,k})} \quad (5)$$

$$c_i = \sum_{j=1}^{|x|} \alpha_{i,j} h_j \quad (6)$$

The attention component of the decoder (termed as non-monotonic attention) in this architecture attends to all the encoder states $h_j : 1 \leq j \leq |x|$ for every decoding time step $1 \leq i \leq |y|$. The model is trained by minimizing the negative log likelihood between predicted and reference translation tokens from the training set.

2) *Monotonic Attention and Policy Network:* [6]–[10] proposed different types of monotonic attention. In terms of architecture, the monotonic attention used in AliBaStr-MT is similar to the non-monotonic attention except that it is equipped with a read/write module, which indicates to the model either read further input or attend to the available encoder states and generate translation. In monotonic attention, the read/write module models the Bernoulli distribution.

$$e_{i,j} = \text{MonotonicEnergy}(s_{i-1}, h_j) \quad (7)$$

$$p_{i,j} = \text{Sigmoid}(e_{i,j}) \quad (8)$$

$$z_{i,j} \approx \text{Bernoulli}(p_{i,j}) \quad (9)$$

If $z_{i,j} = 0$, the model reads the input and updates the encoder state, j is incremented while i and decoder states remain unchanged; if $z_{i,j} = 1$, the model writes the translation, updates the decoder states and increments the decoder step i in addition to updating encoder states and encoder step j . The training of monotonic attention involves the following

- Generate the $p_{i,j}$ for all $1 \leq i \leq |y|$ and $1 \leq j \leq |x|$.
- Compute the expected alignment using the $p_{i,j}$ ($\alpha_{i,j}$ and $\beta_{i,j}$ as presented in [8]).

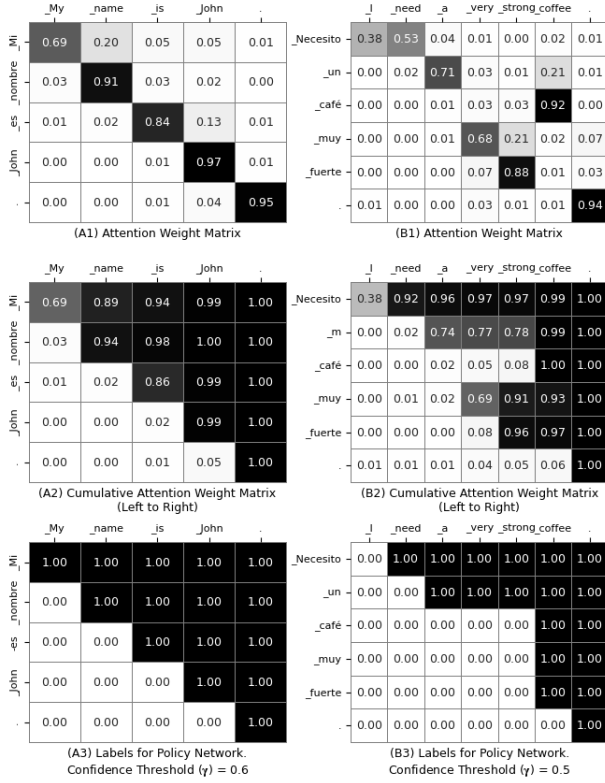


Fig. 3. Conversion of Attention weight matrix to Policy label matrix for supervised training of policy network.

- Compute the context vector c_i .
- Output translation using c_i .

Here, the $p_{i,j}$ are computed first followed by the attention and output projection. In our approach presented in this paper, we follow the standard seq-to-seq training (eq. 1-6). Our proposed model, as shown in Figure 2(c), is also similar to the non-streaming model of Figure 2(a) except that the decoder has a read/write module making the non-monotonic attention behave as the monotonic attention. The attention module only looks at the input tokens consumed so far. It does not have access to the future input tokens (using a mask as described in section III-B3). The policy network is simply a binary classifier over the encoder and decoder states that predicts the probability based on the similarity between the embeddings.

3) *Read/Write Policy Learning*: The architecture of the AliBaStr-MT read/write module is similar to [6], [7] except that the read/write module is not part of the decoder layer. The read/write module is added on top of the decoder after all the decoder layers. We compute the current decoder state at step i and compute the energy function on the current decoder state.

$$e'_{i,j} = \text{MonotonicEnergy}(s_i, h_j) \quad (10)$$

$$p'_{i,j} = \text{Sigmoid}(e'_{i,j}) \quad (11)$$

$$z'_{i,j} \approx \text{Bernoulli}(p'_{i,j}) \quad (12)$$

For training, the labels are generated by the Policy Label Generator (Figure 2 (b)) which uses the attention scores from

the pretrained non-streaming model. If we are able to find the point in the attention weight matrix (with some probability threshold < 1) where all the input tokens have been consumed that are required to generate the i_{th} output token, we can construct a label matrix for direct training of the policy network. For example, consider the input English sentence in Figure 3(A) which is translated in Spanish (using the pretrained non-streaming MT model). If we consider the attention weight matrix of this sentence as shown in the following Figure 3(a1), we can observe that the soft alignment points for each output token can be converted to a policy label matrix using a cumulative attention threshold (Figure 3(a2)). The policy label matrix (Figure 3(a3)) can then be used in the supervised training of the policy network. This example shows a 1:1 mapping of input tokens and output tokens. Figures 3(b) show a more complex example.

4) *Cumulative Attention Confidence Threshold (γ)*: To convert the cumulative attention weight matrix into a discrete policy label matrix, we use an attention confidence threshold (γ). γ ranges from 0 to 1. $\gamma=1$ means a non-streaming case where the policy label generator does not generate a "write" label until it has seen a full sentence. Any value $\gamma < 1$ enables generation of translation with fewer input tokens. The "1" in the label matrix at position (i, j) indicates a point where the model can generate the y_i output token after reading $x_{1..j}$ input tokens. If $\alpha_{i,j}$ represents the attention weights in the attention matrix A , the label matrix L has the following form

$$F = \text{cumsum}(A, \text{dim} = 1) \quad (13)$$

$$L = \text{genlabel}(F) \quad (14)$$

where,

$$l_{i,j} = \begin{cases} 1 & f_{i,j} \geq \gamma \\ 0 & \text{Otherwise} \end{cases} \quad (15)$$

The policy label generator makes sure that the label matrix is monotonic in terms of input/output by enforcing the following.

$$l_{i,j} = \begin{cases} 0 & l_{i-1,j} = 0 \\ l_{i,j} & \text{Otherwise} \end{cases} \quad (16)$$

During training, the output y_i is not generated with fewer input tokens than y_{i-1} . For the example in Figure 3(A), $\gamma=0.6$ was used to generate the policy label matrix. However, $\gamma=0.6$ is pretty specific to this example. The example in Figure 3(B) shows $\gamma = 0.5$ is a good choice.

Furthermore, there is a possibility that the highest $\alpha_{i,j}$ is not covered by the confidence threshold (γ) given that $\gamma < 1$. The Policy Label Generator makes sure that the *write* action is not generated until the highest $\alpha_{i,j}$ has been observed. This is achieved by adding a constraint to eq. 15 as follows

$$l_{i,j} = \begin{cases} 1 & f_{i,j} \geq \gamma \ \& \ \arg \max_{j'} (\alpha_{i,j'}) \leq j \\ 0 & \text{Otherwise} \end{cases} \quad (17)$$

The read/write module operates as a binary classifier, relying on a calibration threshold (δ) to convert probability scores into

clear read/write decisions. Using δ instead of the training-time parameter γ , we can dynamically adjust the quality-latency trade-off during inference, making it easier to identify the optimal model configuration.

C. ASR and MT Integration

To achieve low-latency and high-quality translation, it is essential to carefully integrate the ASR with AliBaStr-MT. Both ASR and MT models in our architecture are designed for real-time processing, allowing them to handle streaming data. The ASR model employs beam search to generate transcripts, which can introduce latency due to its tendency to delay output until a hypothesis is finalized. To mitigate this issue, we have implemented a design that translates both PARTIAL and FINAL hypotheses generated by the beam search.

The beam search produces the top-performing hypothesis as PARTIAL at regular intervals, while FINAL hypothesis is generated when all beams converge on a common prefix. We utilize the PARTIAL hypotheses solely for displaying translations on the screen, whereas FINAL hypotheses are used for playback during conversations. To prevent excessive latency, we have introduced a mechanism that force-finalizes the beams if they fail to produce a FINAL hypothesis within a predetermined time frame (1.5 seconds in our experiments).

We also employ a beam search for AliBaStr-MT. Unlike previous studies on SMT ([6]–[10]), which relied solely on greedy decoding, our approach incorporates streaming beam search decoding. This implementation is based on [32]. Similarly to the ASR, we have introduced force finalization logic into the MT beam search to reduce latency due to finalization delay. Specifically, if hypotheses are not finalized after a predetermined number of output tokens, the beams are automatically finalized.

In our streaming translation design, punctuation prediction is an integral component of the RNN-T model. The AliBaStr-MT model can process input word-by-word, generating translations based on the read/write module. In particular, it resets its internal states when it predicts the end of a sentence, denoted by the punctuation (!?). In contrast, the non-streaming baseline model requires a complete input sentence to produce accurate translations. If only a partial sentence or short phrase is provided, the quality of the translation is compromised. When integrated with ASR, the non-streaming MT model translates the FINAL hypothesis only when it corresponds to a full sentence. However, it can also translate PARTIAL hypotheses, which may consist of shorter fragments. Since these PARTIAL translations are intended for display purposes only, some compromise in quality is acceptable for PARTIALs.

IV. EXPERIMENTS

We evaluate the performance of our speech-to-text translation approach directly on-device on English-Spanish, English-French, English-Italian directions. These language pairs offer more simultaneous translation opportunities than any other language pair. We evaluate the models’ quality/runtime and present the results for both directions for each pair, i.e.

English→X and X→English. For evaluation, we choose internal real-life conversation (RealConv) and publicly available FLEURS [33] datasets.

A. ASR Model Training

We follow a similar approach to [29] for training our ASR system. We use an internal dataset of de-identified video data. The dataset consists of 40k hours of audio. For training, we normalize the training transcripts i.e., all the transcripts are lower-cased, numbers/acronyms are verbalized and all punctuations are removed except for (‘,’ ‘?’ ‘!’ ‘.’). We train a bilingual ASR model for each language pair for on-device usage specifically. The ASR model is not only capable of transcribing the audio but can also label the speaker(self/other), can recognize partially overlapping speech, and perform side talk supersession.

B. MT Model Training

For training of the MT models, we use publicly available datasets, e.g. NLLB, Opensubtitle, CCMatrix, Wikipedia [34]–[37] etc. as well as some internal datasets. Our MT models are for on-device usage. To save memory and other resources on-device, the model for each language pair is bilingual and multitasking, serving different translation use cases. The model translates into one of two target languages; translates from ASR output where formatting is missing compared to regular text; supports code switching between languages; acts as a cross-lingual as well as same language Punc/Cap/ITN module.

Bilingual Modeling: To train bilingual models, a target language *identity* token is added as a prefix to input/output to request the desired language in output. The model is then trained using a mix of English→X and X→English examples in the training set.

Surface/ASR Form: Typically, MT models are trained on written text that is properly formatted with punctuation, capitalization, and numerals, known as Surface Form. However, the output from ASR systems, referred to as ASR Form, often lacks formatting, with missing punctuation, lowercase letters, and verbalized numerals. To improve the integration of our MT model with the ASR, we apply a text normalizer to a portion of the training set, converting it into an ASR-like format on the source side. By combining this modified dataset with the original Surface Form training set, we not only improve the model’s ability to translate from ASR output but also introduce cross-lingual capabilities for punctuation, capitalization, and ITN (Inverse Text Normalization).

Copy Through: We also augment a portion of the training set to perform an identity or a copy-through operation, i.e. English→English or X→X. This is to allow the model to perform code-switching when confronted with mixed language on the input. If the source side of the copy-through dataset is also text normalized, we get the same language Punc/Cap/ITN capabilities useful for the code-switching scenario.

Our primary baseline is a cascade ST system with non-streaming full sentence MT which is an Encoder/Decoder model with monotonic encoder and non-monotonic decoder

TABLE I
WORD ERROR RATE (WER) AND SPEAKER ATTRIBUTION ERROR (SA)
OF THE ASR SYSTEM ON ENGLISH-X PAIRS.

Language	Real-life Conversation			FLEURS	
	↓WER (X)	↓WER (EN)	↓SA	↓WER (X)	↓WER (EN)
EN-ES	11.99	7.55	0.01	10.44	4.95
EN-FR	16.61	7.73	0.07	13.90	10.21
EN-IT	15.99	8.08	0.05	10.56	5.89

having 102M parameters. However, for comparison with AliBaStr-MT, we chose a couple other state-of-the-art SMT models. **Wait-k**: Encoder/Decoder model with monotonic encoder and decoder with wait- k policy [21]. The model has around 115M parameters. **EMMA**: Efficient implementation of MMA approach [6], [7] having 124M parameters.

All of the models have the same configuration except for the architecture specific differences like presence of the tiny read/write module etc. The encoder is based on transformer architecture having 12-layers. The decoder is also transformer-based and has 5-layers. The AliBaStr-MT model has around 103M parameters. The model is similar to the non-streaming model except that it has a small read/write module. For AliBaStr-MT, we set $\gamma = 0.5$ for policy label generation and adjust δ during inference to control quality-latency trade-off. For evaluation, we use our streaming beam-search implementation as introduced in section III-C. The non-streaming models use the standard beam search. For real-time streaming translation, the ASR output (FINAL/PARTIAL) is immediately sent to the SMT model for translation except for the last token which is held for the next call. This is done due to the difference in the ASR and MT sentencepiece vocabulary. Finally, we measure the quality of the model using BLEU, and the latency of the streaming models using Average Lag (AL) [21]. We also introduce User Perceived Latency (UPL), which measures the time it takes for the system to produce the first or last output token after receiving the initial or final input audio frame. As its name implies, UPL is directly observable by users, who can easily notice delays at the beginning and end of the sentence. However, detecting latency during the intermediate stages of translation is more challenging. AL presented in the results section is non-computational-aware while UPL compensates for the computational-aware AL.

C. Results

Table I provides an overview of our ASR system’s performance, highlighting the Word Error Rate (WER) for each language pair (English-X). The results are presented for two distinct datasets: RealConv and FLEURS. Notably, the table also includes Speaker Attribution (SA) errors for the RealConv dataset, which measures the accuracy of predicting speaker tags in multi-party speech recognition scenarios.

Table II presents the end-to-end speech-to-text translation results using a beam size of 1. For each pair of languages, the table reports BLEU for both directions together with AL and UPL in seconds. The scores presented for AL and UPL are the average for both directions. AL differs from UPL in the

TABLE II
BLEU AND LATENCY (P50) METRICS OF ST MODELS FOR ENGLISH \leftrightarrow X
DIRECTIONS ON REAL-LIFE CONVERSATION DATASET.

X	Model	↑BLEU		↓Latency		
		X→EN	EN→X	AL	UPL	
					First	Last
ES	Non-Streaming	45.56	46.60	4.39	4.81	1.99
	Wait-4	42.37	40.77	3.34	3.46	2.18
	EMMA	40.84	41.71	3.20	3.19	2.32
	AliBaStr-MT	43.89	45.30	2.74	3.16	1.36
FR	Non-Streaming	43.02	46.32	4.15	4.32	2.07
	Wait-4	40.82	42.46	3.28	3.40	2.46
	EMMA	40.13	43.06	3.34	3.03	2.67
	AliBaStr-MT	42.13	45.06	2.83	3.18	1.67
IT	Non-Streaming	41.70	32.66	4.49	4.86	1.86
	Wait-4	41.18	31.15	3.36	3.58	2.06
	EMMA	38.56	30.85	3.60	3.40	2.44
	AliBaStr-MT	40.99	32.41	2.89	3.43	1.47

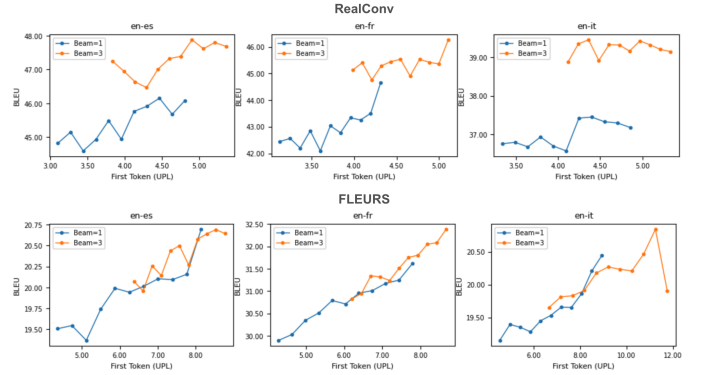


Fig. 4. First Token (UPL) vs. BLEU on RealConv and FLEURS datasets.

sense that AL accounts for the amount of speech (in seconds) taken to generate the output token, while UPL measures the wall-clock time taken to predict the first (last) output token since the system consumed the first (last) input audio frame.

Table II compares the performance of four different models. Our target baseline is, however, the non-streaming model which is the best in quality and has the BLEU that is an upper bound for all the streaming models. In contrast, the latency of the non-streaming model is the worst. On the other hand, the latency of AliBaStr-MT is best across all models. Both non-computational-aware AL and UPL show significant gains compared to the non-streaming model. AliBaStr-MT also has better translation quality (measured in terms of BLEU) than the Wait- k and EMMA models and approaches the quality of the non-streaming model.

Furthermore, Figure 4 presents the latency-quality trade-off curve using the AliBaStr-MT model. These trend curves are generated using $\delta = (0.50, 0.55, 0.60, \dots, 1.0)$. For each pair of languages, the graph presents the trend at beam sizes of 1 and 3. For beam size > 1 , we kept the force finalization token limit to 5 i.e., whenever a beam has 5 or more output tokens, the beams are finalized and only the highest scoring beam that has the maximum number of output tokens is retained. As expected, there is a slight difference in BLEU at beam sizes of 1 and 3. However, with increasing beam size, the model

latency also increases (note that this difference is not captured by the non computation-aware AL which only measures the amount of speech taken to generate output token). In addition to the RealConv dataset, Figure 4 also presents the trend for the FLEURS dataset. For FLEURS dataset, the trend is similar to the RealConv dataset, as UPL increases, the BLEU increases. The highest UPL and BLEU is at $\delta = 1.0$ where the AliBaStr-MT model performs similarly to the non-streaming model.

REFERENCES

- [1] N. Arivazhagan, C. Cherry, I. Te, W. Macherey, P. Baljekar, and G. Foster, “Re-translation strategies for long form, simultaneous, spoken language translation,” in *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 7919–7923.
- [2] D. Li, T. I. N. Arivazhagan, C. Cherry, and D. Padfield, “Sentence boundary augmentation for neural machine translation robustness,” in *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021, pp. 7553–7557.
- [3] P. Koehn and R. Knowles, “Six challenges for neural machine translation,” in *Proceedings of the First Workshop on Neural Machine Translation*, T. Luong, A. Birch, G. Neubig, and A. Finch, Eds. Vancouver: Association for Computational Linguistics, Aug. 2017, pp. 28–39. [Online]. Available: <https://aclanthology.org/W17-3204/>
- [4] K. Cho, B. van Merriënboer, D. Bahdanau, and Y. Bengio, “On the properties of neural machine translation: Encoder-decoder approaches,” in *Proceedings of SSST-8, Eighth Workshop on Syntax, Semantics and Structure in Statistical Translation*, D. Wu, M. Carpuat, X. Carreras, and E. M. Vecchi, Eds. Doha, Qatar: Association for Computational Linguistics, Oct. 2014, pp. 103–111. [Online]. Available: <https://aclanthology.org/W14-4012/>
- [5] J. Pouget-Abadie, D. Bahdanau, B. van Merriënboer, K. Cho, and Y. Bengio, “Overcoming the curse of sentence length for neural machine translation using automatic segmentation,” in *Proceedings of SSST-8, Eighth Workshop on Syntax, Semantics and Structure in Statistical Translation*, D. Wu, M. Carpuat, X. Carreras, and E. M. Vecchi, Eds. Doha, Qatar: Association for Computational Linguistics, Oct. 2014, pp. 78–85. [Online]. Available: <https://aclanthology.org/W14-4009/>
- [6] X. Ma, A. Sun, S. Ouyang, H. Inaguma, and P. Tomasello, “Efficient monotonic multihead attention,” *arXiv preprint arXiv:2312.04515*, 2023.
- [7] X. Ma, J. M. Pino, J. Cross, L. Puzon, and J. Gu, “Monotonic multihead attention,” 2019.
- [8] N. Arivazhagan, C. Cherry, W. Macherey, C.-C. Chiu, S. Yavuz, R. Pang, W. Li, and C. Raffel, “Monotonic infinite lookback attention for simultaneous machine translation,” in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 2019, pp. 1313–1323.
- [9] C.-C. Chiu and C. Raffel, “Monotonic chunkwise attention,” in *International Conference on Learning Representations*, 2018.
- [10] C. Raffel, M.-T. Luong, P. J. Liu, R. J. Weiss, and D. Eck, “Online and linear-time attention by enforcing monotonic alignments,” in *International conference on machine learning*. PMLR, 2017, pp. 2837–2846.
- [11] N. Moritz, R. Xie, Y. Gaur, K. Li, S. M. Zeeshan Ahmed, F. Seide, and C. Fuegen, “Transcribing and translating, fast and slow: Joint speech translation and recognition,” in *ICASSP 2025 - 2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2025, pp. 1–5.
- [12] R. Dabre and H. Song, “NICT’s cascaded end-to-end speech translation systems using whisper and IndicTrans2 for the Indic task,” in *Proceedings of the 21st International Conference on Spoken Language Translation (IWSLT 2024)*, E. Salesky, M. Federico, and M. Carpuat, Eds. Bangkok, Thailand (in-person and online): Association for Computational Linguistics, Aug. 2024, pp. 17–22. [Online]. Available: <https://aclanthology.org/2024.iwslt-1.3/>
- [13] S. Communication, L. Barrault, Y.-A. Chung, M. C. Meglioli, D. Dale, N. Dong, M. Duppenhaler, P.-A. Duquenne, B. Ellis, H. Elsahar, J. Haasheim, J. Hoffman, M.-J. Hwang, H. Inaguma, C. Klaiber, I. Kulikov, P. Li, D. Licht, J. Maillard, R. Mavlyutov, A. Rakotoarison, K. R. Sadagopan, A. Ramakrishnan, T. Tran, G. Wenzek, Y. Yang, E. Ye, I. Evtimov, P. Fernandez, C. Gao, P. Hansanti, E. Kalbassi, A. Kallet, A. Kozhevnikov, G. M. Gonzalez, R. S. Roman, C. Touret, C. Wong, C. Wood, B. Yu, P. Andrews, C. Balioglu, P.-J. Chen, M. R. Costa-jussà, M. Elbayad, H. Gong, F. Guzmán, K. Heffernan, S. Jain, J. Kao, A. Lee, X. Ma, A. Mourachko, B. Peloquin, J. Pino, S. Popuri, C. Ropers, S. Saleem, H. Schwenk, A. Sun, P. Tomasello, C. Wang, J. Wang, S. Wang, and M. Williamson, “Seamless: Multilingual expressive and streaming speech translation,” 2023. [Online]. Available: <https://arxiv.org/abs/2312.05187>
- [14] S. Zhang, Q. Fang, S. Guo, Z. Ma, M. Zhang, and Y. Feng, “StreamSpeech: Simultaneous speech-to-speech translation with multi-task learning,” in *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, L.-W. Ku, A. Martins, and V. Srikumar, Eds. Bangkok, Thailand: Association for Computational Linguistics, Aug. 2024, pp. 8964–8986. [Online]. Available: <https://aclanthology.org/2024.acl-long.485/>
- [15] J. Iranzo-Sánchez, J. Jorge, P. Baquero-Arnal, J. A. Silvestre-Cerdà, A. Giménez, J. Civera, A. Sanchis, and A. Juan, “Streaming cascade-based speech translation leveraged by a direct segmentation model,” *Neural Networks*, vol. 142, pp. 303–315, 2021.
- [16] P. Bahar, T. Bieschke, and H. Ney, “A comparative study on end-to-end speech to text translation,” in *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2019, pp. 792–799.
- [17] J. Iranzo-Sánchez, J. Iranzo-Sánchez, A. Giménez, J. Civera, and A. Juan, “Segmentation-free streaming machine translation,” *Transactions of the Association for Computational Linguistics*, vol. 12, pp. 1104–1121, 2024. [Online]. Available: <https://aclanthology.org/2024.tacl-1.61/>
- [18] R. Fukuda, K. Sudoh, and S. Nakamura, “Speech segmentation optimization using segmented bilingual speech corpus for end-to-end speech translation,” *arXiv preprint arXiv:2203.15479*, 2022.
- [19] D. Soboleva, O. Skopek, M. Sajjalik, V. Cărbune, F. Weissenberger, J. Proskurnia, B. Prisacari, D. Valcarce, J. Lu, R. Prabhavalkar, and B. Miklos, “Replacing human audio with synthetic audio for on-device unspoken punctuation prediction,” in *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021, pp. 7653–7657.
- [20] F. Dalvi, H. Sajjad, S. Vogel, and N. Durrani, “Incremental decoding and training methods for simultaneous translation in neural machine translation,” in *Proceedings of NAACL-HLT*, 2018, pp. 493–499.
- [21] M. Ma, L. Huang, H. Xiong, R. Zheng, K. Liu, B. Zheng, C. Zhang, Z. He, H. Liu, X. Li *et al.*, “Stacl: Simultaneous translation with implicit anticipation and controllable latency using prefix-to-prefix framework,” in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 2019, pp. 3025–3036.
- [22] S. Guo, S. Zhang, and Y. Feng, “Glancing future for simultaneous machine translation,” in *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2024, pp. 11 386–11 390.
- [23] K. Cho and M. Esipova, “Can neural machine translation do simultaneous translation?” *arXiv preprint arXiv:1606.02012*, 2016.
- [24] A. Grissom II, H. He, J. Boyd-Graber, J. Morgan, and H. Daumé III, “Don’t until the final verb wait: Reinforcement learning for simultaneous machine translation,” in *Proceedings of the 2014 Conference on empirical methods in natural language processing (EMNLP)*, 2014, pp. 1342–1352.
- [25] J. Gu, G. Neubig, K. Cho, and V. O. Li, “Learning to translate in real-time with neural machine translation,” in *15th Conference of the European Chapter of the Association for Computational Linguistics, EACL 2017*. Association for Computational Linguistics (ACL), 2017, pp. 1053–1062.
- [26] O. Press and N. A. Smith, “You may not need attention,” *arXiv preprint arXiv:1810.13409*, 2018.
- [27] A. Alinejad, H. S. Shavarani, and A. Sarkar, “Translation-based supervision for policy generation in simultaneous neural machine translation,” in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, M.-F. Moens, X. Huang, L. Specia, and S. W.-t. Yih, Eds. Online and Punta Cana, Dominican Republic: Association for Computational Linguistics, Nov. 2021, pp. 1734–1744. [Online]. Available: <https://aclanthology.org/2021.emnlp-main.130/>
- [28] S. Papi, M. Negri, and M. Turchi, “Attention as a guide for simultaneous speech translation,” in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, A. Rogers, J. Boyd-Graber, and N. Okazaki, Eds. Toronto, Canada: Association for Computational Linguistics, Jul. 2023, pp. 13 340–13 356. [Online]. Available: <https://aclanthology.org/2023.acl-long.745/>

- [29] J. Lin, N. Moritz, R. Xie, K. Kalgaonkar, C. Fuegen, and F. Seide, "Directional speech recognition for speaker disambiguation and cross-talk suppression," in *Proc. Interspeech*, vol. 2023, 2023, pp. 3522–3526.
- [30] J. Lin, N. Moritz, Y. Huang, R. Xie, M. Sun, C. Fuegen, and F. Seide, "Agadir: Towards array-geometry agnostic directional speech recognition," in *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2024, pp. 11 951–11 955.
- [31] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [32] R. Rabatin, F. Seide, and E. Chang, "Navigating the minefield of mt beam search in cascaded streaming speech translation," *arXiv preprint arXiv:2407.11010*, 2024.
- [33] A. Conneau, M. Ma, S. Khanuja, Y. Zhang, V. Axelrod, S. Dalmia, J. Riesa, C. Rivera, and A. Bapna, "Fleurs: Few-shot learning evaluation of universal representations of speech," in *2022 IEEE Spoken Language Technology Workshop (SLT)*. IEEE, 2023, pp. 798–805.
- [34] H. Schwenk, G. Wenzek, S. Edunov, E. Grave, and A. Joulin, "Ccmatrix: Mining billions of high-quality parallel sentences on the web," *arXiv preprint arXiv:1911.04944*, 2019.
- [35] A. Fan, S. Bhosale, H. Schwenk, Z. Ma, A. El-Kishky, S. Goyal, M. Baines, O. Celebi, G. Wenzek, V. Chaudhary *et al.*, "Beyond english-centric multilingual machine translation," *Journal of Machine Learning Research*, vol. 22, no. 107, pp. 1–48, 2021.
- [36] P. Lison and J. Tiedemann, "Opensubtitles2016: Extracting large parallel corpora from movie and tv subtitles," 2016.
- [37] K. Wolk and K. Marasek, "Building subject-aligned comparable corpora and mining it for truly parallel sentence pairs," *Procedia Technology*, vol. 18, pp. 126–132, 2014.