

REINA: Regularized Entropy Information-Based Loss for Efficient Simultaneous Speech Translation

Anonymous ACL submission

Abstract

Simultaneous Speech Translation (SimulST) systems stream in audio while simultaneously emitting translated text or speech. Such systems face the significant challenge of balancing translation quality and latency. We introduce a strategy to optimize this tradeoff: wait for more input only if you gain information by doing so. Based on this strategy, we present Regularized Entropy INformation Adaptation (REINA), a novel loss to train an adaptive policy using an existing non-streaming translation model. We derive REINA from information theory principles and show that REINA helps push the reported Pareto frontier of the latency/quality tradeoff over prior works. Utilizing REINA, we train a SimulST model on French, Spanish and German, both from and into English. Training on only open source or synthetically generated data, we achieve state-of-the-art (SOTA) streaming results for models of comparable size. We also introduce a metric for streaming efficiency, quantitatively showing REINA improves the latency/quality trade-off by as much as 21% compared to prior approaches, normalized against non-streaming baseline BLEU scores.

1 Introduction

Simultaneous Speech Translation (SimulST) involves real-time translation of speech in one language into text in another. This extends the simpler speech-to-text-translation (S2TT) task, which involves translation with the full context of an entire speech clip. While S2TT allows for offline applications, conversational environments such as voice or video chat necessitate SimulST models to facilitate real-time communication across language barriers.

Recently, End-to-end (E2E) S2TT models have largely superseded traditional cascaded approaches, which link separate Automatic Speech Recognition (ASR) and Machine Translation (MT) systems. E2E models mitigate error propagation and reduce

latency by directly mapping source speech to target text (Communication et al., 2023; Peng et al., 2024; Radford et al., 2023; Puvvada et al., 2024).

Rather than training SimulST models from scratch, most works take advantage of advances in S2TT research by adapting a non-streaming S2TT model into a SimulST model (Communication et al., 2023; Chen et al., 2024; Zhang et al., 2024; Papi et al., 2024). To transition from non-streaming S2TT to SimulST introduces the challenge of balancing translation quality and latency. This requires a policy to decide whether to wait for more input (READ) or generate output (WRITE) (Gu et al., 2017). This problem is particularly difficult in the translation setting because different languages can have different word orderings, requiring differing amounts of context before a suitable translation can occur. Many approaches have been formulated to determine this READ/WRITE policy, from baking the policy into the model architecture itself via monotonic attention mechanisms (Ariavazhagan et al., 2019; Ma et al., 2020c; Communication et al., 2023), or having a separate module to dictate the policy (Chen et al., 2024). However, these existing approaches suffer from issues including poor translation quality compared to non-streaming models and expensive, numerically unstable training (Communication et al., 2023; Chen et al., 2024).

In this paper, we address the problem of efficiently training high quality SimulST models. The major contributions of this paper can be summarized as follows:

New Policy Training Technique. We propose Regularized Entropy INformation Adaptation (REINA), a new technique for policy training that can efficiently convert non-streaming Speech-to-Text Translation (S2TT) models into simultaneous S2TT (SimulST) models. REINA is guided by an approximation of mutual information derived from the S2TT model’s log probabilities on partial

versus full audio and is shown to produce higher quality policies than existing methods.

Trained with open-source data. We train an E2E S2TT model with REINA on 130k hours of open-source data. Based on empirical studies, this model achieves SOTA streaming translation performance.

Streaming efficiency evaluation metric. We propose a new evaluation metric to better compare SimulST models. This metric normalizes the streaming translation quality against the performance of the underlying non-streaming model, allowing for a fairer assessment of the capabilities of the streaming policy itself.

2 Related Work

While learning adaptive policies for SimulST is a fairly new research area, it builds on a rich body of non-streaming S2TT work. In this section, we outline the S2TT foundation for SimulST and then move into a discussion of SimulST policies.

Training Speech to Text Translation Models The literature around S2TT contains many large-scale, powerful models including Whisper (Radford et al., 2023), SeamlessM4T (Communication et al., 2023), Canary (Puvvada et al., 2024), and the Open Whisper-style Speech Model (OWSM) (Peng et al., 2024). These models vary in architecture (e.g., Whisper’s Transformer, OWSMv3.1’s E-Branchformer) and training data scale, ranging from Canary’s 86k hours (leveraging pseudo-labels) and OWSM’s 180k hours of public data to Whisper’s 680k hours of web data and SeamlessM4T’s 600k hours of synthetically aligned data.

Due to the relative scarcity of parallel ST data compared to ASR or MT corpora, multi-task learning (MTL) is widely adopted (Ye et al., 2022; Chen et al., 2024; Communication et al., 2023). Auxiliary tasks like ASR and neural machine translation (NMT) are jointly trained with S2TT to improve representations and leverage abundant text or speech data. Furthermore, contrastive learning techniques, such as in ConST (Ye et al., 2022), are used to explicitly bridge the modality gap between speech and text representations by encouraging similarity between corresponding speech segments and their transcriptions. There is a notable gap in the literature between industry work leveraging massive proprietary datasets and less resourced research making heavy use of MTL to get the most

out of smaller data scales. As OWSM (Peng et al., 2024) bridges this gap for the non-streaming setting, we aim to do the same for SimulST. We are one of the first SimulST works to leverage large-scale open source data that we train on with an MTL framework including MT and ASR tasks.

Streaming Policy Learning Transitioning S2TT models to SimulST introduces the challenge of learning a READ/WRITE policy that balances translation quality and latency. Fixed policies like wait-k are simple to implement but are usually sub-optimal due to the mismatch between the sampling rate of the input audio frames and the frequency of outputted words (Ma et al., 2020b).

On the other hand, adaptive policies dynamically adjust decisions based on context. Prior works have integrated the policy within the model architecture, such as Transducer models (Graves, 2012; Xue et al., 2022), which inherently support streaming via monotonic alignment, or models using monotonic attention mechanisms like MMA (Ariavazhagan et al., 2019) or EMMA (Communication et al., 2023). Monotonic alignment methods afford greater expressivity, but they tend to be excessively expensive to compute at train time and suffer from both poor numerical stability and difficulty in converging. For more details on the difficulties of training monotonic attention methods see appendix D.

Other SimulST works avoid complex, explicit policies, instead generating aligned data with which to directly train SimulST models (Labiausse et al., 2025; Fu et al., 2025; Deng et al., 2025). These works often use existing models such as NMT models (Labiausse et al., 2025) or LLM’s (Fu et al., 2025; Deng et al., 2025) as teachers to create synthetically aligned data for streaming training. Such models can afford simpler architectures without policy networks. That said, SimulST models deriving their policies from generated data are often limited in their streaming performance based on the quality of the teacher model.

Explicit policy training Other adaptive strategies decouple the policy from the translation model (Chen et al., 2024; Gu et al., 2017; Zhang et al., 2024). Some leverage signals from pre-trained offline models, such as using reinforcement learning (RL) to directly optimize the quality-latency trade-off (Gu et al., 2017). Although these methods simplify the learning problem by decoupling the policy from the translation model, they require explicit supervision from a suitable metric and are often suboptimal. RL is hard to stabi-

lize and efficiently train, especially in cases like SimulST, with further modifications required for stabilizing the policy head training and no guarantee of convergence (Gu et al., 2017).

Closely related to our work is the divergence-guided approach of DiG-SST (Chen et al., 2024). DiG-SST trains a lightweight policy module using the expected divergence between output distributions conditioned on partial versus complete input, estimated from a non-streaming S2TT model. This approach is efficient to train and directly optimizes for the SimulST task. Nevertheless, DiG-SST’s formulation fails to make use of valuable information from ground truth labels when computing divergence scores. In REINA, we propose an improved formulation of a similar concept, yielding better streaming results.

3 Model

In this section, we introduce REINASStream, a low latency SimulST model trained on large-scale, open source data. We present our architecture and loss functions visually in figure 1.

3.1 Policy Learning

To learn an effective READ/WRITE policy, we introduce a new loss function: Regularized Entropy INformation Adaptation (REINA). REINA enables us to adapt a non-streaming speech translation model into a streaming, SimulST model with minimal extra training.

First, we outline the problem more formally. Suppose that we are translating an input audio stream a into a target language with a streaming chunk size of j frames. Given a partial audio recording a_t at frame t and previously emitted tokens s_1, s_2, \dots, s_n , we need to decide whether to produce token s_{n+1} (WRITE) or wait for another audio chunk (READ). If we READ, we consume another frame of audio, giving us a_{t+1} , whereas if we WRITE, we gain a token, yielding the same audio a_t but tokens s_1, \dots, s_n, s_{n+1} . Our policy must learn to make READ/WRITE decisions that maximize translation quality while minimizing the latency with which we emit each token. This recursive setup appears to lend itself to a dynamic programming type of optimization, as in Seamless (Communication et al., 2023) or the Transducer architecture (Graves, 2012). However, in practice, optimizing over all possible READ/WRITE sequences results in expensive, numerically unsta-

ble training.

Instead, we start from a core idea: we should wait for more audio (i.e. READ) if and only if we gain information by doing so. We formalize this notion using mutual information theory. Given audio a of length T and ground truth translation token sequence $S = (s_1, \dots, s_N)$, after writing $n < N$ tokens and listening to $t < T$ timesteps of audio, we can express the information gained about the next token s_{n+1} by waiting for the rest of the input audio as

$$\mathcal{F}(a, S, n, t) := I(s_{n+1}; a_T, S_n) - I(s_{n+1}; a_t, S_n) \quad (1)$$

where I is the symbol for mutual information. We can then construct an ideal READ/WRITE policy π_α on top of this quantity: $\pi_\alpha(a, S, n, t)$ returns READ when $\mathcal{F}(a, S, n, t) > \alpha$ and WRITE otherwise. We can then adjust α to control the latency quality tradeoff. This policy READs exactly when the information gained exceeds a given threshold. We can rewrite $\mathcal{F}(a, S, n, t)$ using mutual information equations (Barber and Agakov, 2004) as follows

$$\begin{aligned} \mathcal{F}(a, S, n, t) &:= I(s_{n+1}; a_T, S_n) - I(s_{n+1}; a_t, S_n) \\ &= H(s_{n+1}) - H(s_{n+1}|a_T, S_n) \\ &\quad - [H(s_{n+1}) - H(s_{n+1}|a_t, S_n)] \\ &= H(s_{n+1}|a_t, S_n) - H(s_{n+1}|a_T, S_n) \\ &= \mathbb{E} [\log p(s_{n+1}|a_T, S_n) - \log p(s_{n+1}|a_t, S_n)] \end{aligned} \quad (2)$$

Although we might not have access to $\log p(s_{n+1}|a_T, S_n)$ or $\log p(s_{n+1}|a_t, S_n)$, we are able to estimate these via random sampling of the log-probabilities of the base S2TT model as $\log \hat{p}(s_{n+1}|a_T, S_n)$ and $\log \hat{p}(s_{n+1}|a_t, S_n)$, which are the log-probabilities of the next label token computed when passing the full audio and partial audio through the model respectively. Note that these are also the negatives of the cross-entropy losses obtained when running the model on the full and partial audio.

We now have an estimate $\hat{\mathcal{F}}(a, S, n, t)$. However, this estimate cannot be computed during inference time as it requires access to the full target text, so we instead formulate a heuristic to estimate the information gain. We train a model parameterized by θ to estimate a heuristic $q_\theta = q(a, S, n, t|\theta)$ that yields policy $\hat{\pi}_\alpha$, which returns READ if and only if $q_\theta > \alpha$. We train q_θ to strongly correlate with $\mathcal{F}(a, S, n, t)$ by maximizing the covariance

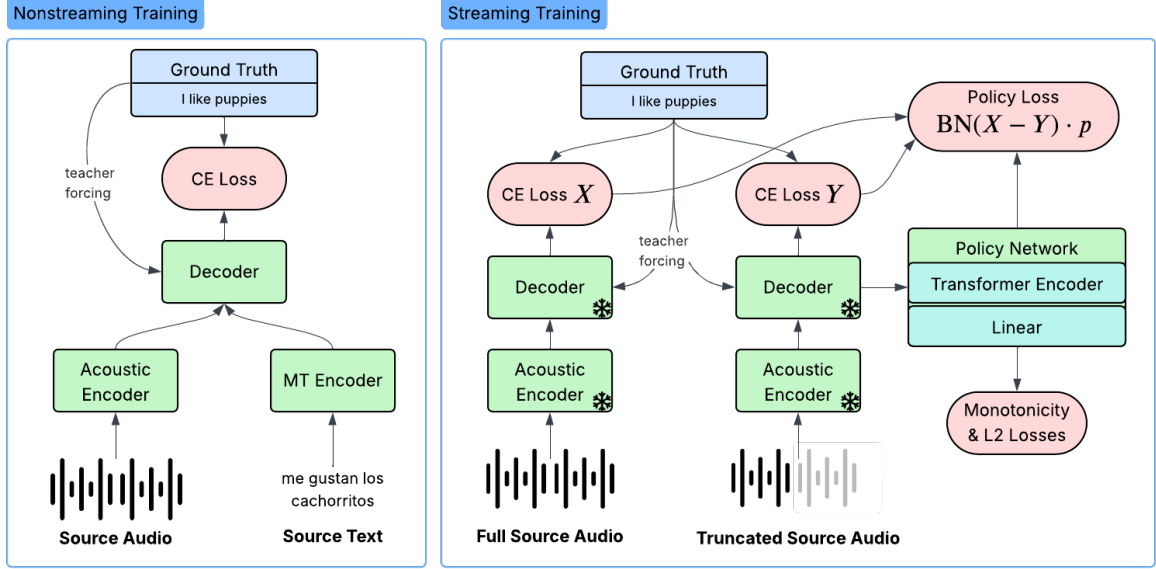


Figure 1: Non-streaming and streaming training procedures for REINASStream. For non-streaming training we use a trainable MT encoder to train on parallel NMT data. During streaming training we a) pass a full audio and truncated audio through the model, b) compute the cross-entropy (CE) loss of each, c) predict a policy using the policy network on top of the partial-audio output of the decoder, and finally d) calculate the REINA loss using the CE terms and policy predictions.

between q_θ and the estimate $\hat{\mathcal{F}}(a, S, n, t)$. This gives us the following optimization problem:

$$\begin{aligned} & \max_{\theta} \left(\text{Cov}(q_\theta, \hat{\mathcal{F}}(a, S, n, t)) \right) \\ & = \max_{\theta} \left(\mathbb{E} \left[q_\theta \cdot \hat{\mathcal{F}}(a, S, n, t) \right] - \mathbb{E}[q_\theta] \cdot \mathbb{E} \left[\hat{\mathcal{F}}(a, S, n, t) \right] \right) \end{aligned} \quad (3)$$

We can simplify this expression by normalizing the information gain estimate $\hat{\mathcal{F}}(a, S, n, t)$ to have zero mean, making $\mathbb{E}[q_\theta] \cdot \mathbb{E}[\hat{\mathcal{F}}(a, S, n, t)]$ evaluate to 0. We achieve this at train time by normalizing over each batch, resulting in the final optimization problem:

$$\begin{aligned} & \max_{\theta} \left[\mathbb{E} \left[q_\theta \cdot \text{BN} \left(\hat{\mathcal{F}}(a, S, n, t) \right) \right] \right] \\ & = \min_{\theta} \left[q_\theta \cdot \text{BN} \left[\log \hat{p}_t^{s_{n+1}} - \log \hat{p}_T^{s_{n+1}} \right] \right] \end{aligned} \quad (4)$$

where BN is the batchnorm operator, and we apply a shorthand $\log \hat{p}_T^{s_{n+1}} = \log \hat{p}(s_{n+1}|a_T, S_n)$ and $\log \hat{p}_t^{s_{n+1}} = \log \hat{p}(s_{n+1}|a_t, S_n)$ for brevity.

The loss that optimizes the above optimization problem is

$$\mathcal{L}_p = \frac{1}{N} \sum_{n=0}^{N-1} q_\theta^n \cdot \text{BN} \left[\log \hat{p}_t^{s_{n+1}} - \log \hat{p}_T^{s_{n+1}} \right] \quad (5)$$

where $q_\theta^n = q_\theta(a, S, n, t|\theta)$. This loss maximizes the covariance between our estimate of information gain $\hat{\mathcal{F}}$ and our heuristic estimator network q_θ . This is the "Policy Loss" in figure 1. Since $\log \hat{p}_t(s_{n+1})$ and $\log \hat{p}_T(s_{n+1})$ are the negative cross-entropy losses from the decoder given partial and full audio respectively, in the diagram we write \mathcal{L}_p as a difference of cross-entropy loss terms.

To ensure stable training and a reasonable learned policy, we add monotonicity and L_2 regularization terms, both of which are outlined below. First, we note that at inference time, after predicting a READ, we predict no further tokens. Therefore, any WRITES after a READ are semantically meaningless.

To better align train and inference, we add a weak monotonicity constraint on q_θ to encourage the probability of READ to increase uniformly across each token sequence S :

$$\mathcal{L}_m = \frac{1}{N} \sum_{n=1}^N \left[\max \left(\max_{m < n} \{ q_\theta^m \} - q_\theta^n - \epsilon, 0 \right) \right] \quad (6)$$

The loss function \mathcal{L}_m serves as a regularization term to shape the learned heuristic q_θ used in the policy $\hat{\pi}_\alpha(a, S, n, t)$. It encourages the sequence of q_θ values for target text tokens, ordered by index n , to be approximately non-decreasing

($q_\theta^n \geq \max_{m < n} \{q_\theta^m\} - \epsilon$), an inductive bias distinct from merely estimating information gain. This imposed monotonicity biases the policy towards a “commitment” behavior, encouraging it to commit to stop predicting new tokens and READ new audio instead once the threshold is crossed for an earlier token.

Finally, we add a simple L_2 regularization penalty: $\mathcal{L}_r = \frac{1}{N} \sum_{n=1}^N (q_\theta^n)^2$. We find this is required to prevent q_θ values from exploding to infinity during training. Putting all the terms together, we get the full REINA loss: $\mathcal{L}_{\text{REINA}} = \mathcal{L}_p + \mathcal{L}_m + \lambda \mathcal{L}_r$. In our work, we set $\lambda = 0.05$, but find that final model performance is not very sensitive to changes in λ . We will cover how we train with $\mathcal{L}_{\text{REINA}}$ in section 3.4.

3.2 Non-Streaming Architecture

Next, we describe the architecture of our base non-streaming S2TT model, which is also outlined in figure 1. We employ a straightforward architecture comprising an acoustic encoder and a text decoder. At train-time, we also use an extra text encoder to facilitate a MT training task.

For the acoustic encoder, we adopt Whisper Medium (Radford et al., 2023) and do not freeze its weights during training. The randomly initialized transformer decoder performs cross-attention over the acoustic encoder’s final layer hidden states and predicts text tokens. We adopt the generic multilingual tokenizer from Mistral 7B (Jiang et al., 2023) but learn our own embedding dictionary. We augment the vocabulary with language ids such as <en> or <fr> so we can direct the decoder to predict tokens of a specific target language by prefixing the token sequence with the language id. We apply a learned positional encoding similar to Time2Vec (Kazemi et al., 2019) on the acoustic encoder outputs to give the decoder a notion of sequence ordering. The decoder is trained with a cross-entropy loss.

At train-time, to support MT loss calculation, we add a randomly initialized, trainable T5 text encoder (Raffel et al., 2020). We pass source-language text through and then have the decoder cross-attend to the T5 last layer hidden states while predicting target language text. This facilitates a machine translation task designed to improve the quality of the decoder by making use of paired MT data.

Whisper Medium contains 307M parameters, the text decoder has 101M, and the MT encoder has

38M for a total of 445M trainable parameters at train time and 408M at inference-time. While this is a larger parameter count than academic works like Dig-SST (Chen et al., 2024), Stream-Speech (Zhang et al., 2024), or Stream-Att (Papi et al., 2024), it is still much smaller than most industry systems like SeamlessM4T (Communication et al., 2023) or Hibiki (Labiausse et al., 2025). While higher parameter counts demonstrably increase translation quality, they also make it difficult to train and deploy models to large numbers of users in the wild. We target the middle ground between small and large systems in literature, yielding a model with both the translation quality and computational efficiency to be usable in real-world chat settings.

3.3 Streaming Architecture

We augment the non-streaming model with a policy network that makes binary READ/WRITE decisions for each output token. This is the architecture component that REINA serves to train. The policy network is a small transformer encoder applied to the last layer hidden states from the decoder. We apply a single linear layer on top of the transformer with output dimension 1 and sigmoid activation in order to make the binary READ/WRITE decisions. We apply a causal attention mask on the policy network at both train and inference time. The policy network contains only 6M trainable parameters, bringing the total parameter count for the streaming model at train and inference time to 414M.

3.4 Training

We train REINASStream in 3 stages: 1) Learn non-streaming S2TT 2) Adapt to truncated audios 3) Learn a streaming policy.

In the first stage, we train on several tasks at once in order to effectively leverage available data to train the speech translation model. All data samples contain some subset of the following information: source language l_i^s , source language audio a_i , source language transcription T_i^s , target language l_i^t , and target language transcription T_i^t . Using this data, we have three training tasks.

ASR For samples with (a_i, l_i^s, T_i^s) , we pass a_i through the acoustic encoder, then decode to tokens in source language l_i^s . We compute cross-entropy loss using label transcript T_i^s yielding loss \mathcal{L}_{asr} .

NMT For samples with (T_i^s, l_i^s, T_i^t) , we pass the source transcript through the T5 text encoder and then decode into target language l_i^t . We com-

pute the cross-entropy loss using label transcript T_i^t yielding loss \mathcal{L}_{nmt} .

S2TT For samples with (a_i, l_i^t, T_i^t) , we pass a_i through the acoustic encoder, then decode into target language l_i^t . We compute cross-entropy loss using label transcript T_i^t yielding loss \mathcal{L}_{s2tt} .

For the first stage training, we mix data supporting all tasks into every batch and minimize a sum of all the losses: $\mathcal{L} = \mathcal{L}_{asr} + \mathcal{L}_{nmt} + \mathcal{L}_{s2tt}$.

Before training the policy network on top of the base model, we add a second step to ensure quality estimation of $\log \hat{p}(s_{n+1}|S_n, a_t)$ with partial audios a_t . In this phase, we fine-tune the speech translation model on randomly truncated audios using the same loss function \mathcal{L} .

Lastly, we train the policy network by minimizing \mathcal{L}_{REINA} and freezing all other parameters. We only train on S2TT data samples as our goal is to learn a policy best for streaming speech translation rather than streaming ASR or MT.

3.5 Data

We seek to bridge the gap between models trained on large-scale, proprietary datasets and those trained on small-scale open-source data. We leverage a variety of publicly available data sources plus synthetic data generation to produce a large-scale training set.

In this iteration, we focus only on en→de, fr, es and de, fr, es→en language directions because of their data availability. In future work, we plan to expand to lower resourced languages. For audio datasets, we draw from Multilingual LibriSpeech (MLS) (Pratap et al., 2020), Mosel (Gaido et al., 2024), CVSS-C (Jia et al., 2022), MUST-C (Di Gangi et al., 2019). We list details on these datasets in table 1. We further augment the MLS dataset by translating its transcripts using an in-house NMT model to produce S2TT data to train on. We also augment our dataset with text-to-text MT training data from CCMatrix (Schwenk et al., 2021). We use 10M samples per language pair from CCMatrix for a total of 60M samples.

3.6 Inference Policy

Once our model is trained, we employ a streaming beam search to perform inference. We split input audios into 0.25s chunks and inference the model on all audio up to the current chunk in sequence.

We first pick a policy threshold α to control the quality-latency tradeoff while streaming. Each iteration of the search, we run the policy network on

all beams and kill those with predictions less than α . If the total number of waited beams exceeds the beam size times a patience factor, or all beams are waited at once, we end the search. After ending the search, we return the killed beam with the highest average log probability.

Once we have reached the end of the input audio, we stop using the policy network and beam search until we hit the EOS token, using the same patience factor logic to decide when to end the search.

4 Experiments

In this section, we detail how we ran our experiments and present our results compared with existing work.

4.1 Experimental Setup

To train REINASTream, we follow the procedure outlined in section 3.4, starting by training our non-streaming S2TT model on a mixture of all training datasets. We train on 24 A100-80G GPUs for 5 days using an AdamW optimizer (Loshchilov and Hutter, 2019) with hyperparameters given in appendix C. We use the exact same setup to train the second truncated stage, but only train for 2 days. Lastly, we perform the REINA training stage with nearly the same setup, except this time using an inverse square root learning rate scheduler with 5k warmup steps. This final training runs 20 epochs in under 12 hours.

We also re-implement DiG-SST’s divergence-based loss from (Chen et al., 2024) based on the description in the paper and train the policy network with that loss using the same configuration as we train REINA.

We evaluate on language pairs {fr, de, es} → en on CVSS-C, and on language pairs en → {fr, de, es} on MUST-C. We use a beam size of 3 with no length penalty, a streaming chunk size of 0.25s, and a patience factor of 3. We sweep across several thresholds for the policy network in order to measure our model’s tradeoff between latency and accuracy. Choosing the right thresholds to obtain an informative sweep is done purely through trial and error (more details in appendix C).

We compare to existing works on each dataset that we believe to be at or near to state of the art at the time of writing. On MUST-C, we compare to Dig-SST (Chen et al., 2024), the work in literature closest to ours, and another strong SimulST competitor called DiSeg (Zhang and Feng, 2023).

Split	Dataset	Source Language: en				Source Language: de		Source Language: es		Source Language: fr	
		Target: de	Target: es	Target: fr	Target: en	Target: en	Target: de	Target: en	Target: es	Target: en	Target: fr
Train	MLS	13,789	13,785	13,787	41360	1,637	1,637	713	713	984	984
	Must-C	386	476	468	1330	-	-	-	-	-	-
	CVSS-C	-	-	-	-	184	184	113	113	264	264
	MOSEL	-	-	-	19,245	-	22,804	-	19,373	-	22,835
Dev	Must-C	2.47	2.49	2.49	-	-	-	-	-	-	-
	CVSS-C	-	-	-	-	21	-	22	-	22	-
Test	Must-C	4	4	4	-	-	-	-	-	-	-
	CVSS-C	-	-	-	-	22	-	23	-	23	-

Table 1: Consolidated Dataset Hours (Source Audio) by Split and Dataset, Grouped by Source and Target Language.

On CVSS-C, we compare to StreamSpeech (Zhang et al., 2024), a recent work showing strong simultaneous S2ST performance on CVSS with a small model. As StreamSpeech only reports ASR-BLEU whereas we report text BLEU, we are unable to make a fair comparison. That said, StreamSpeech is a stronger system than most in the literature on CVSS-C, so we believe the results are still worth showing. In appendix E, we also show some results on the FLEURS (Conneau et al., 2023) dataset with significant caveats. For all comparisons, we use self-reported results from the original papers.

We also perform several ablations on the REINA training stage of REINASTream to demonstrate the utility of the different parts of the REINA loss and make fairer comparisons. We train four model variants:

REINA Our standard training procedure including all S2TT training datasets with the standard REINA loss function.

REINA w/o monotonicity Just like **REINA** but without the monotonicity term in the loss. Trained on all S2TT datasets.

REINA (MUST-C only) We train the policy network with the full REINA loss on only the MUST-C dataset for a fairer comparison to Dig-SST.

Dig-SST (Our impl. MUST-C only) We train our own implementation of the Dig-SST loss on top of the non-streaming REINASTream model on only the MUST-C dataset.

4.2 Evaluation Metrics

We measure translation accuracy using BLEU as implemented in the SacreBLEU package (Post, 2018) and Average Lag (AL), which we implement ourselves based on the original paper (Ma et al., 2020a). Most existing works plot AL vs BLEU curves by interpolating between several (AL, BLEU) points generated by evaluating their model with different streaming settings (Chen et al., 2024) (Papi et al., 2024) (Zhang et al., 2024). We include such graphs for all of our evaluations in appendix B.

We contend that this evaluation does not adequately disentangle a model’s non-streaming translation quality from its streaming ability. We observe that many comparisons in the literature pitch models with superior non-streaming BLEU as being better at streaming due to having a higher BLEU vs AL curve, when in reality the difference may be accounted for entirely by the non-streaming BLEU difference. Such comparisons do little to show the quality of the actual READ/WRITE policy.

We wish to show that REINA adapts S2TT models of any non-streaming performance to SimulST models with minimal degradation in translation quality. To this end, we introduce a new metric, Normalized Streaming Efficiency (NoSE), to measure streaming performance across the entire quality/latency spectrum normalized by non-streaming translation quality. As shown in the diagram in appendix A, we measure the area under the AL/BLEU curve, bounded on the left and right by $x < y$ respectively, divided by the area under the non-streaming BLEU line. It is important to note that while we require the x and y bounds to ensure the metric is well-defined, NoSE is heavily dependent on them. For our analysis, we simply pick the smallest x and largest y for which our work and the works we compare to all have reported values, yielding the widest possible range for which all models have a defined AL/BLEU curve. We recommend that future works using NoSE do report the bounds used in their calculations.

4.3 Results

We present our NoSe scores in table 2 and selected operating points in table 3 for MUST-C and CVSS results respectively, with full AL-BLEU tradeoff curves in appendix B.

Unfortunately, CVSS-C is not commonly evaluated against in SimulST literature, so we are unable to compare to many other works. However, with an average utterance length of 4.9 seconds, CVSS is the only dataset out of the three comprising pri-

marily shorter audios, which are quite common in conversational SimulST use-cases. This makes it an important benchmark for SimulST systems.

MUST-C			
Model	en→de	en→fr	en→es
Bounds $[x, y]$	[1.193, 2.629]	[1.11, 2.887]	[1.144, 2.611]
Dig-SST (Original)	0.923	0.945	0.927
DiSeg	0.881	-	0.880
Dig-SST (Our impl. MUST-C only)	0.737	0.846	0.740
REINA (MUST-C only)	0.957	0.972	0.980
REINA	0.945	0.966	0.976

CVSS			
Model	de→en	fr→en	es→en
Bounds $[x, y]$	[1.955, 5.039]	[1.637, 5.169]	[1.806, 5.587]
StreamSpeech*	0.842	0.886	0.837
REINA w/o monotonicity	0.976	0.980	0.982
REINA	0.974	0.983	0.981

Table 2: NoSE (\uparrow) values for the MUST-C and CVSS-C datasets. Bounds $[x, y]$ for NoSE are specified per language pair. Best values per language pair are bolded. Results for other works are taken from their corresponding papers. *Note StreamSpeech only reports ASR-BLEU, so values for StreamSpeech represent ASR-BLEU scores rather than BLEU.

4.3.1 Quantitative results

On both MUST-C and CVSS-C, we observe REINA outperforms all competing methods on all language splits at low latencies. Significantly, this holds for the "MUST-C only" model with policy network trained only on MUST-C, as well as the REINA model trained on all datasets, demonstrating our streaming performance gains are not merely attributable to increased data scale when training the policy, but come from the improved objective. The only exception to this, as seen in table 3, is German, where DigSST is slightly better at higher latencies than REINA. Still, this shows REINA excels at lower latency streaming, even when its non-streaming BLEU is lower than competitors.

The MUST-C-only REINA model yields NoSE scores 3.0% higher than Dig-SST and 8.9% higher than DiSeg, highlighting the advantages of REINA’s mutual information formulation. Our implementation of DiG-SST performs far below every other model in evals, suggesting there may be details to the authors’ implementation that we failed to recreate.

4.3.2 Ablations on Monotonicity Loss

In figure 2, we observe that the REINA model outperforms the REINA w/o monotonicity model exclusively at low latencies, indicating that monotonicity is useful on the most aggressive streaming settings. For example, at about 35 BLEU, AL de-

Lang.	REINA (musc only)		DIGSST (reported)	
	AL \downarrow	BLEU \uparrow	AL \downarrow	BLEU \uparrow
De	1.01	21.44	1.08	21.13
	1.59	23.71	1.45	23.25
	2.24	24.32	1.83	24.29
Es	0.86	26.92	0.90	23.92
	1.16	29.68	1.27	26.74
	1.51	30.38	1.66	27.90
Fr	0.77	33.13	1.11	30.51
	1.24	36.29	1.26	32.61
	1.88	37.73	2.00	35.65

Table 3: Comparison of REINA (musc only) with DIGSST (reported) on En→{De, Es, Fr} speech translation. Selected operating points are reported.

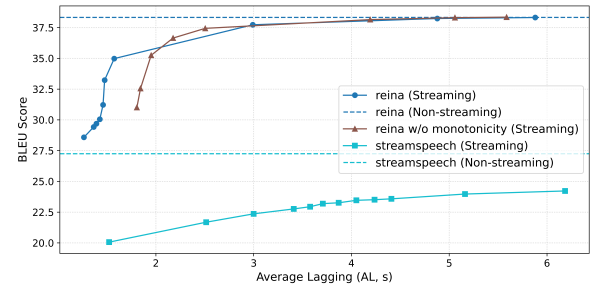


Figure 2: AL/BLEU curve on es-en split of the CVSS-C dataset.

creases from 1.95 to 1.57, a 19% improvement. We hypothesize this is because monotonicity forces the policy to decide on a clear boundary of when to READ when the information gain waffles between timesteps. Ultimately, we find the monotonicity loss improves low-latency streaming with a negligible impact to overall streaming efficiency.

5 Conclusion and Future Work

In this paper, we present a new method for SimulST. We introduce the REINA loss function that enables cheaply converting non-streaming speech translation models into streaming ones. We conduct extensive experiments over several datasets, showing REINA outperforms the state of the art in SimulST conversion. We also propose a new metric, NoSE, to improve the state of evaluation of SimulST systems. Ultimately, we train a large-scale system entirely on open source or synthetically generated data, encouraging further research into scaling SimulST.

The next step to enable real-time crosslingual interaction is to extend REINASTream into a simultaneous speech to speech translation (SimulS2ST) model. This is achievable by using a high quality, low latency, streaming text-to-speech model as a synthesizer. We are presently working on extending REINASTream to the SimulS2ST use-case.

6 Limitations

6.1 Evaluation Datasets

Evaluations in literature are split between several datasets, namely MUST-C and CVSS-C. These datasets alone already present a variety of issues—for example, MUST-C is no longer publicly available for download and many of its audios do not contain any speech at all. Furthermore, none of these datasets contain large-scale training data, resulting in a variety of data augmentation and training strategies across the literature, making it difficult to fairly compare models trained on the same data. We train on a custom mix of synthetic and human-labeled data that is not directly comparable to any one prior piece of research. We evaluate on all three of the aforementioned datasets, but have no guarantees that our data cleaning or inference procedures are identical to existing work.

6.2 Implementation of Evaluation

On the implementation side, we also likely do not achieve implementation-parity with existing work. Most SimulST research is built on top of the fairseq package, obfuscating implementation details. Fairseq performs a number of optimizations under the hood that are not fully described in literature but affect model quality and evaluation results. Examples include checkpoint averaging on the training side and cleaning audios that are too short or long on the evaluation side. To improve reproducibility, control over our architecture, and eventual ease-of-serving, we take the route of implementing REINASStream, our entire evaluation suite, and our data preprocessing from scratch. This means there are likely implementation details that differ between our evaluation and works built on top of fairseq. In the future, where possible, we should aim to run open source existing works through our evaluation pipeline.

6.3 Translation Quality Metrics

We only use BLEU to report translation quality, which is widely known to not correlate well with human perception of translation quality. While BLEU allows us to compare to existing work, we plan to move toward more robust neural metrics such as COMET or LLM-based metrics.

6.4 Synthetic vs Human Labeled Data

Our training data for ASR, MT, and S2TT tasks is largely synthetically generated by existing ASR

and NMT models. This bounds REINASStream’s performance based on the quality of the models used to aid in the generation of our datasets. Seamless tackles this data scarcity issue by using their SONAR alignment model to mine pairs from real-world data. We may try out the same technique in future work. Another mitigation strategy could be to make heavier use of human-labeled NMT and ASR data that is available online. We may also be able to use of RLHF methods becoming popular in the NMT world to improve our translation quality with less human-labeled data.

6.5 Low Resource Languages

We train and evaluate only on four languages all with relatively similar grammatical structure and substantial open source data to train on. It is widely known that low training data resources heavily affect translation quality. Similarly, translating between different language families is much harder than translating between related languages. This is especially true in the SimulST case where different word orders make the READ/WRITE policy’s task more difficult and necessitate higher latency for the same translation quality. For instance, Japanese has a subject-object-verb structure whereas English uses subject-verb-object. When translating simultaneously from Japanese to English, a good policy likely will have to learn to wait until late in the Japanese utterance to determine the verb and use it in the English translation. In future work, we will extend our training and evaluation to multiple language families and include low resource language pairs.

6.6 REINA Mutual Information Approximation

In the REINA loss, we only approximate $\log p(s_{n+1}|S_n, a_t) - \log p(s_{n+1}|S_n, a_T)$ using our translation model. However, there is an easy way we could compute the true value rather than merely approximating. If we generated synthetic labels S_N by sampling from our translation model and then trained the policy using the REINA loss on these synthetic labels, $\log \hat{p}(s_{n+1}|S_n, a_t) - \log \hat{p}(s_{n+1}|S_n, a_T) = \log p(s_{n+1}|S_n, a_t) - \log p(s_{n+1}|S_n, a_T)$ would hold. We did not have enough time to experiment with this idea before submitting and plan on pursuing it in future work.

6.7 Setting Policy Threshold

We find that determining the appropriate threshold α for our policy $\hat{\pi}_\alpha$ at inference time is rather difficult. Depending on the training data and evaluation dataset, the range of thresholds that produce reasonable AL/BLEU tradeoff varies. This forces us to try out many different thresholds at inference time, which wastes time and compute. In the future, we plan on investigating ways of cheaply automatically finding usable policy thresholds.

6.8 FLEURS Evaluation

Along with MUST-C and CVSS, FLEURS is an S2TT dataset evaluated on by SimulST works, most notably Seamless. We originally sought to include FLEURS in our evaluation suite, but ran into several challenges.

First, as FLEURS contains multiple speakers speaking each utterance, we were unsure to what level we should dedupe in order to compare to other work, and could not find reliable documentation specifying this. Second, we found that especially when translating $\text{en} \rightarrow \text{X}$, we had poor non-streaming BLEU scores, suggesting some domain mismatch between our other datasets and FLEURS. Third, we found Seamless difficult to compare against because a) they modify the test dataset using VAD, b) they report a very small range of AL scores to compute NoSE over, c) Due to scarce documentation, we were unable to reproduce Seamless’ results locally. Due to these difficulties, we leave our caveat-ridden FLEURS evaluation in appendix E.

References

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.

David Barber and Felix Agakov. 2004. The im algorithm: a variational approach to information maximization. *Advances in neural information processing systems*, 16(320):201.

Xinjie Chen, Kai Fan, Wei Luo, Linlin Zhang, Libo Zhao, Xinggao Liu, and Zhongqiang Huang. 2024. [Divergence-guided simultaneous speech translation](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(16):17799–17807.

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

Alexis Conneau, Min Ma, Simran Khanuja, Yu Zhang, Vera Axelrod, Siddharth Dalmia, Jason Riesa, Clara Rivera, and Ankur Bapna. 2023. [Fleurs: Few-shot learning evaluation of universal representations of speech](#). In *2022 IEEE Spoken Language Technology Workshop (SLT)*, pages 798–805.

Keqi Deng, Wenxi Chen, Xie Chen, and Philip C. Woodland. 2025. [Simul2s-llm: Unlocking simultaneous inference of speech llms for speech-to-speech translation](#). *Preprint*, arXiv:2504.15509.

Mattia A. Di Gangi, Roldano Cattoni, Luisa Bentivogli, Matteo Negri, and Marco Turchi. 2019. [MuST-C: a Multilingual Speech Translation Corpus](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2012–2017, Minneapolis, Minnesota. Association for Computational Linguistics.

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

Marco Gaido, Sara Papi, Luisa Bentivogli, Alessio Brutti, Mauro Cettolo, Roberto Gretter, Marco Matasoni, Mohamed Nabih, and Matteo Negri. 2024. [MOSEL: 950,000 hours of speech data for open-source speech foundation model training on EU languages](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 13934–13947, Miami, Florida, USA. Association for Computational Linguistics.

Alex Graves. 2012. [Sequence transduction with recurrent neural networks](#). *Preprint*, arXiv:1211.3711.

Jiatao Gu, Graham Neubig, Kyunghyun Cho, and Victor O.K. Li. 2017. [Learning to translate in real-time with neural machine translation](#). In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 1053–1062, Valencia, Spain. Association for Computational Linguistics.

Ye Jia, Michelle Tadmor Ramanovich, Quan Wang, and Heiga Zen. 2022. Cvss corpus and massively multilingual speech-to-speech translation. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 6691–6703.

873	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.	931
874		932
875		933
876		934
877		
878		935
879		936
880		937
881	Seyed Mehran Kazemi, Rishab Goel, Sepehr Eghbali, Janahan Ramanan, Jaspreet Sahota, Sanjay Thakur, Stella Wu, Cathal Smyth, Pascal Poupart, and Marcus A. Brubaker. 2019. Time2vec: Learning a vector representation of time . <i>CoRR</i> , abs/1907.05321.	938
882		939
883		940
884		941
885		
886	Tom Labiausse, Laurent Mazar��, Edouard Grave, Patrick P��rez, Alexandre D��fossez, and Neil Zeghidour. 2025. High-fidelity simultaneous speech-to-speech translation. <i>arXiv preprint arXiv:2502.03382</i> .	942
887		943
888		944
889		945
890	Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization . In <i>International Conference on Learning Representations</i> .	946
891		
892		947
893	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.	948
894		949
895		950
896		951
897		952
898		
899		953
900	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.	954
901		955
902		956
903		957
904		958
905		959
906		960
907		961
908		
909	Xutai Ma, Juan Miguel Pino, James Cross, Liezl Puzon, and Jiatao Gu. 2020c. Monotonic multihead attention . In <i>International Conference on Learning Representations</i> .	962
910		963
911		964
912		965
913	Sara Papi, Marco Gaido, Matteo Negri, and Luisa Bentivogli. 2024. Streamatt: Direct streaming speech-to-text translation with attention-based audio history selection. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3692–3707.	966
914		967
915		968
916		969
917		970
918		971
919	Yifan Peng, Jinchuan Tian, William Chen, Siddhant Arora, Brian Yan, Yui Sudo, Muhammad Shakeel, Kwanghee Choi, Jiatong Shi, Xuankai Chang, and 1 others. 2024. Owsm v3. 1: Better and faster open whisper-style speech models based on e-branchformer. In <i>Proc. Interspeech 2024</i> , pages 352–356.	972
920		973
921		974
922		975
923		976
924		977
925		978
926	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.	979
927		980
928		981
929		982
930		983
	Vineel Pratap, Qiantong Xu, Anuroop Sriram, Gabriel Synnaeve, and Ronan Collobert. 2020. Mls: A large-scale multilingual dataset for speech research. In <i>Proc. Interspeech 2020</i> , pages 2757–2761.	984
		985
		986
	Krishna C Puvvada, Piotr ��elasko, He Huang, Oleksii Hrinchuk, Nithin Rao Koluguri, Kunal Dhawan, Somshubra Majumdar, Elena Rastorgueva, Zhehuai Chen, Vitaly Lavrukhin, and 1 others. 2024. Less is more: Accurate speech recognition & translation without web-scale data. In <i>Proc. Interspeech 2024</i> , pages 3964–3968.	
	Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. Robust speech recognition via large-scale weak supervision. In <i>Proceedings of the 40th International Conference on Machine Learning, ICML’23</i> . JMLR.org.	
	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>Journal of machine learning research</i> , 21(140):1–67.	
	Holger Schwenk, Guillaume Wenzek, Sergey Edunov, Edouard Grave, Armand Joulin, and Angela Fan. 2021. CCMatrix: Mining billions of high-quality parallel sentences on the web . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 6490–6500, Online. Association for Computational Linguistics.	
	Jian Xue, Peidong Wang, Jinyu Li, Matt Post, and Yashesh Gaur. 2022. Large-scale streaming end-to-end speech translation with neural transducers. In <i>Proc. Interspeech 2022</i> , pages 3263–3267.	
	Rong Ye, Mingxuan Wang, and Lei Li. 2022. Cross-modal contrastive learning for speech translation. In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 5099–5113.	
	Shaolei Zhang, Qingkai Fang, Shoutao Guo, Zhengrui Ma, Min Zhang, and Yang Feng. 2024. Stream-speech: Simultaneous speech-to-speech translation with multi-task learning. In <i>Proceedings of the 62th Annual Meeting of the Association for Computational Linguistics (Long Papers)</i> . Association for Computational Linguistics.	
	Shaolei Zhang and Yang Feng. 2023. End-to-end simultaneous speech translation with differentiable segmentation . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , page 7659–7680. Association for Computational Linguistics.	
	A NoSE Score	
	To clarify our definition of the NoSE score, we provide a diagram in figure 3.	

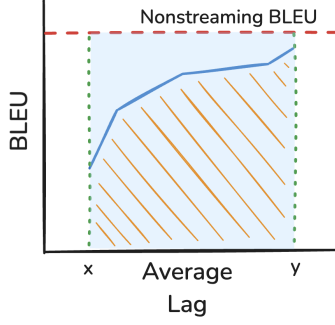


Figure 3: We define NoSE as the area of the orange shaded region divided by the area of the blue rectangle.

B Latency vs BLEU Evaluations

When evaluating our models on each dataset, we sweep across several policy thresholds. We find that the unique features of each dataset necessitate different thresholds to get the best results. In table 4 we show the list of thresholds we used for inferring each model on each dataset. Unfortunately, the only way we have to determine these thresholds is trial and error. We had to perform several offline inference sweeps across thresholds before finding a sweep that a) allowed us to compare to other works b) contained points in the AL/BLEU curve that were representative of how our model might be likely to be used in production.

We show AL vs BLEU graphs for all evaluations that appear in the paper in figures 4 and 5.

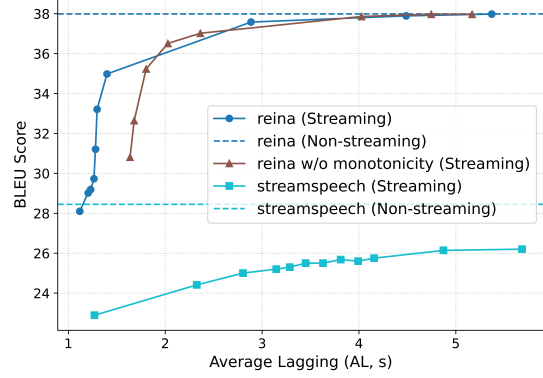
C Hyperparameters

We report model and dataset hyperparameters in table 5.

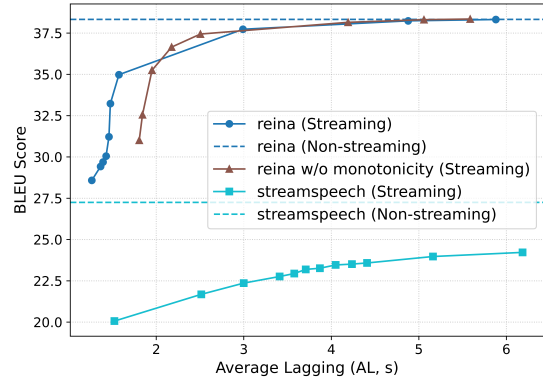
D EMMA Discussion

In our paper, we touch on the complexities of implementing streaming via monotonic attention as in the EMMA method used in (Communication et al., 2023). In this section, we wish to provide some more background on the computational and numerical challenges that arise from training with EMMA.

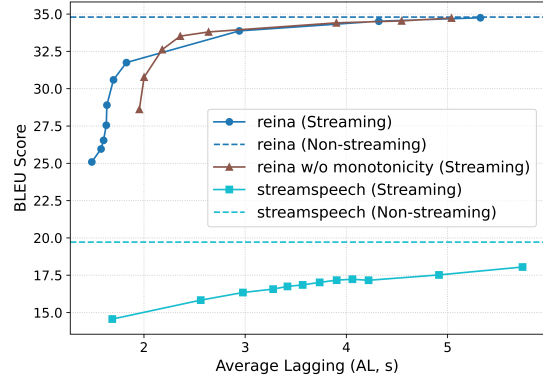
In our research, we implemented EMMA from scratch on top of a non-streaming REINASStream model. The train-time computations for EMMA requires computing a matrix of size $[\text{batch_size} \times \text{attention_heads} \times \text{num_text_tokens} \times \text{audio_sequence_length} \times \text{audio_sequence_length}]$ within each cross-attention layer of the decoder. Seeing as the audio sequence length coming out



(a) French \rightarrow English (fr-en)



(b) Spanish \rightarrow English (es-en)

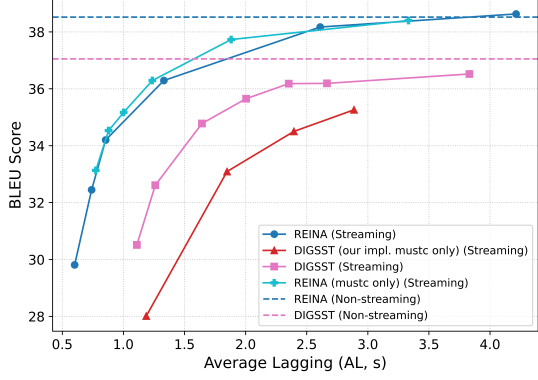


(c) German \rightarrow English (de-en)

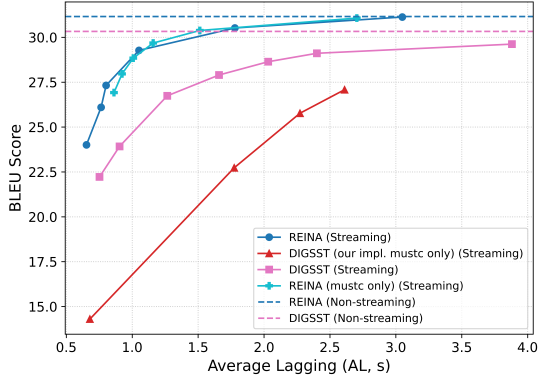
Figure 4: Average Lagging (AL) vs. BLEU score on CVSS-C. Dotted lines represent non-streaming BLEU scores. Note that StreamSpeech only reports ASR-BLEU in their paper, so we report StreamSpeech’s ASR-BLEU rather than BLEU.

Table 4: Thresholds swept for generating AL vs. BLEU points for different models and datasets.

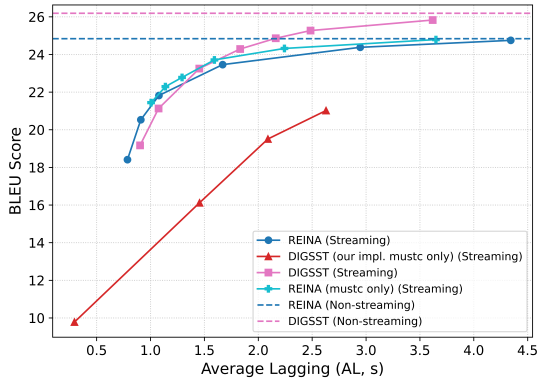
Dataset	Model Name	Thresholds Swept
MUST-C	Dig-SST (Original)	Self-reported results
	DiSeg	Self-reported results
	Dig-SST (Our impl. MUST-C only)	[.475, .5, .51, .52, .53, .54]
	REINA (MUST-C only)	[.935, .94, .9425, .945, .9475, .95]
	REINA	[.97, .975, .976, .977, .978, .979]
CVSS	REINA	[.97, .975, .976, .977, .978, .979, .98, .983, .985, .987]
	REINA w/o monotonicity	[.976, .977, .978, .979, .98, .983, .985, .987]
	StreamSpeech	Self-reported results



(a) English → French (en-fr)



(b) English → Spanish (en-es)



(c) English → German (en-de)

Figure 5: Average Lagging (AL) vs. BLEU score on MUST-C. Dotted lines represent non-streaming BLEU scores.

Table 5: Hyperparameters and Dataset Mixing Ratios used for training the S2TT model and the REINA policy network.

Parameter / Setting	Value
Text Decoder Configuration	
Dimension	512
Attention Heads	8
Number of Layers	16
Feedforward Multiplier	4
Label Smoothing	0.1
Dropout	0.1
Policy Network Configuration (REINA)	
Monotonicity Epsilon (ϵ)*	0.5
L2 Regularization Weight (λ)	0.05
Number of Layers	2
Attention Heads	4
Feedforward Multiplier	4
Dimension	512
Training Configuration	
Learning Rate	0.0001 (fixed)
Weight Decay	0.0001
Gradient Clipping	10.0
Batch Size per Device	16
Gradient Accumulation Steps	2
Effective Batch Size	768
Dataset Mixing: Base S2TT Model Training (Stage 1 & 2)	
Dataset Name	Mixing Ratio
MUST-C (train split)	1
CVSS (train split)	1
MLS (en → X)**	2
MLS (X → en)**	4
CCMatrix	4
Mosel	4
Dataset Mixing: REINA Policy Training (Stage 3)	
Dataset Name	Mixing Ratio
MUST-C (train split)	2
CVSS (train split)	2
MLS (en → X)**	6
MLS (X → en)**	6

* Corresponds to ϵ in the monotonicity loss (Equation 3).

** Derived from Multilingual LibriSpeech (MLS) dataset (see Section 3.5).

of the whisper encoder is 1500, we use 8 attention heads, and a short token sequence may contain around 25 tokens, using fp32 precision, we require 2GB of VRAM for a single cross-attention layer at batch size 1. Ultimately, we were only able to train EMMA with batch size 1 and truncating the encoder output sequence length to 500, despite using A100-80G GPUs. This made for very slow, expensive training.

Furthermore, the EMMA estimation requires computing a cumulative product across the audio sequence length dimension. A cumulative product of 500 small floating point values is numerically unstable and often results in rounding to 0. Perhaps the original authors used a sum of log values instead.

Lastly, as EMMA computes a separate policy for every attention head of every cross-attention layer, it is unclear which one to use for the final inference policy. In the public inference code associated with (Communication et al., 2023), the layer to use for the policy is simply taken as an argument. Attention heads from that layer are aggregated via a max, min, or mean. Empirically, we found that some layers and heads learned useful policies, while the majority did not.

Due to these challenges implementing and training EMMA, we did not include it in our evaluations and pursued simpler streaming methods, motivating us to invent REINA.

We also note that many of the same issues arise with training neural transducer (Graves, 2012) models. We also spent considerable effort implementing transducers for SimulST. As with EMMA, we found transducer methods to be excessively expensive to train and very hard to make converge.

E FLEURS Evaluation

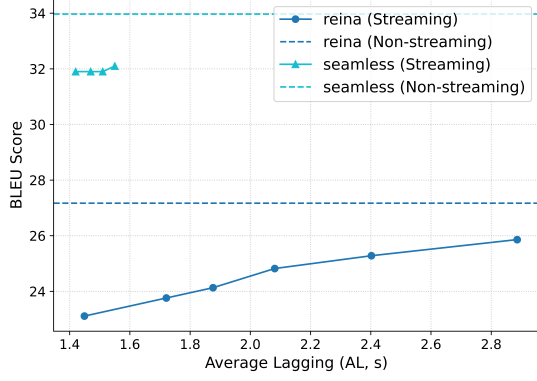
In this section, we present an attempt at comparing to Seamless on FLEURS. We construct the FLEURS dataset for en \rightarrow de, en, fr and de, en, fr \rightarrow en, deduplicating on unique source audio and target language. We trim the dataset using the same Silero VAD Seamless uses. When inferencing, we soon noticed that as our model is trained entirely on audios not trimmed with VAD, it tends to expect an ending silence and over-generates, resulting in decreased BLEU scores. We also attempted inferencing Seamless on our copy of FLEURS using their open source code, but ran into implementation difficulties due to scarce documentation.

In a final attempt at a comparison, we use Seamless’ reported scores from their paper. We use VAD to trim our audios and then augment them with 2 seconds of white noise at the end to help our model. We tested this trick without VAD on our other datasets and found it has a negative impact on BLEU, but we found it was still an improvement over inferencing on VAD-trimmed audios on FLEURS.

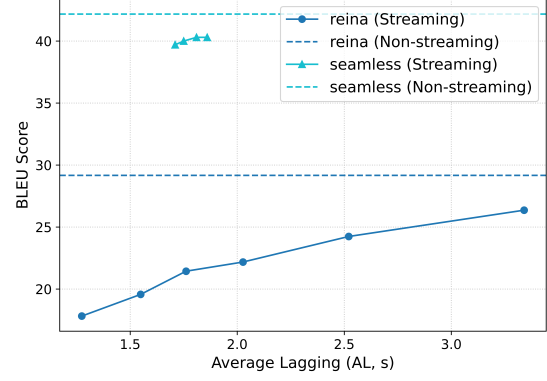
We show NoSE scores in table 6, graphs for the X \rightarrow en direction in figure 6, and graphs for the en \rightarrow X direction in figure 7. As our model is much smaller than Seamless, our non-streaming BLEU is unsurprisingly much lower. Of more interest are the NoSE scores, which show we are comparable to Seamless on the X \rightarrow en directions but worse on en \rightarrow X. We believe this difference is primarily accounted for by the aforementioned VAD issues. We also note that REINA is vastly cheaper and easier to train than Seamless’ EMMA streaming method (see appendix D), meaning that even with comparable streaming quality, REINA is in most cases preferable to EMMA.

Dataset	Item / Model	en→de	en→fr	en→es	de→en	fr→en	es→en
	<i>Bounds</i> $[x, y]$	[1.87, 2.05]	[1.71, 1.86]	[1.83, 1.99]	[1.68, 1.85]	[1.42, 1.55]	[1.36, 1.52]
FLEURS	Seamless	0.914	0.951	0.960	.924	.940	.936
	REINASStream	0.896	0.861	0.916	.943	.866	.936

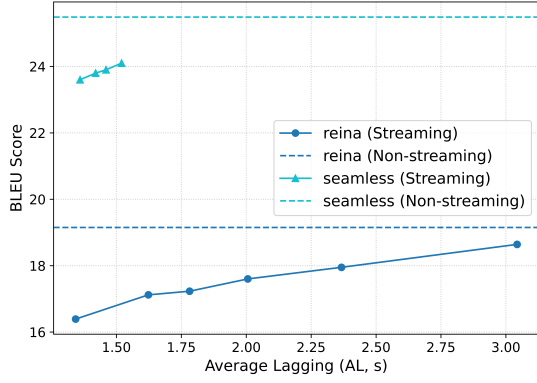
Table 6: Comparison of NoSE (\uparrow) values on FLEURS for Seamless and REINASStream.



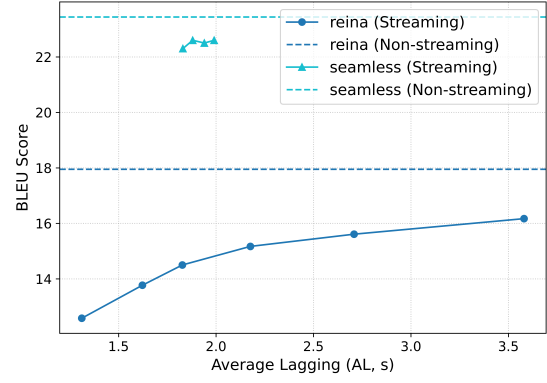
(a) French → English (fr-en)



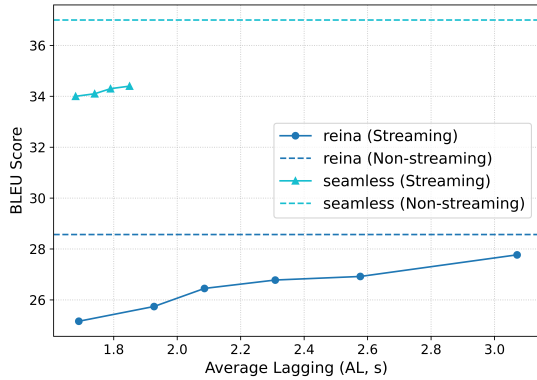
(a) English → French (en-fr)



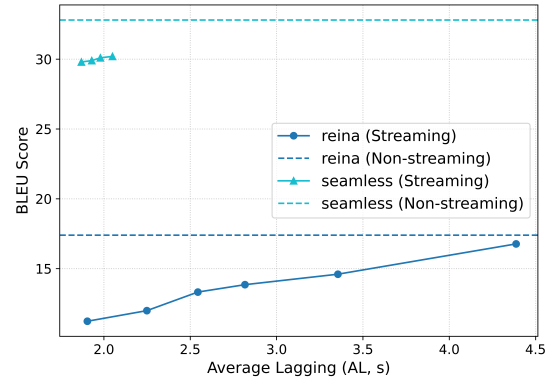
(b) Spanish → English (es-en)



(b) English → Spanish (en-es)



(c) German → English (de-en)



(c) English → German (en-de)

Figure 6: Average Lagging (AL) vs. BLEU score on FLEURS $X \rightarrow en$. Dotted lines represent non-streaming BLEU scores.

Figure 7: Average Lagging (AL) vs. BLEU score on FLEURS $en \rightarrow X$. Dotted lines represent non-streaming BLEU scores.