

LEARNING MONOTONIC ATTENTION IN TRANSDUCER FOR STREAMING GENERATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Streaming generation models are increasingly utilized across various fields, with the Transducer architecture being particularly popular in industrial applications. However, its input-synchronous decoding mechanism presents challenges in tasks requiring non-monotonic alignments, such as simultaneous translation, leading to suboptimal performance in these contexts. In this research, we address this issue by tightly integrating Transducer’s decoding with the history of input stream via a learnable monotonic attention mechanism. Our approach leverages the forward-backward algorithm to infer the posterior probability of alignments between the predictor states and input timestamps, which is then used to estimate the context representations of monotonic attention in training. This allows Transducer models to adaptively adjust the scope of attention based on their predictions, avoiding the need to enumerate the exponentially large alignment space. Extensive experiments demonstrate that our MonoAttn-Transducer significantly enhances the handling of non-monotonic alignments in streaming generation, offering a robust solution for Transducer-based frameworks to tackle more complex streaming generation tasks. Codes are publicly available in supplementary materials.

1 INTRODUCTION

Streaming generation is a widely studied problem in fields such as speech recognition (Raffel et al., 2017; Zhang et al., 2020; Seide et al., 2024), simultaneous translation (Cho & Esipova, 2016; Gu et al., 2017; Seamless Communication et al., 2023), and speech synthesis (Ma et al., 2020a; Zhang et al., 2024; Wang et al., 2024). Unlike modern turn-based large language models, streaming models need to start generating the output before the input is completely read. This necessitates a careful balance between generation quality and latency.

Popular streaming generation methods can be broadly divided into two categories: Attention-based Encoder-Decoder (AED; Bahdanau et al., 2015) and Transducer (Graves, 2012). Streaming AED models adapt the conventional sequence-to-sequence framework (Bahdanau, 2014) to support streaming generation. They often rely on an external policy module to determine the READ/WRITE actions in inference and to direct the scope of cross-attention in training. Examples include Wait- k policy (Ma et al., 2019) and monotonic attention-based methods (Raffel et al., 2017; Arivazhagan et al., 2019; Ma et al., 2020d; 2023a). On the other hand, Transducer models connect the encoder and predictor through a joiner rather than using cross-attention. As shown in Figure 1a, the joiner is designed to synchronize the encoder and predictor by expanding its output vocabulary to include a blank symbol ϵ , which indicates a READ action. Due to the decoupling of the predictor state from the encoder state, READ/WRITE states in Transducer can be represented by a two-dimensional lattice. This allows for the computation of total probabilities using the forward-backward algorithm (Graves, 2012), facilitating end-to-end optimization. Benefited from joint optimization of all potential policies during training, Transducer often demonstrates better performance compared to AED models (Xue et al., 2022; Wang et al., 2023).

During the decoding process of Transducer, each target token is explicitly aligned with a corresponding source token. This input-synchronous decoding property makes the architecture well-suited for tasks like speech recognition, where the input and output align monotonically. However, it poses challenges for non-monotonic alignment tasks such as simultaneous translation (Chuang et al., 2021; Shao & Feng, 2022; Ma et al., 2023c). Due to the decoupled design, Transducer models have limited

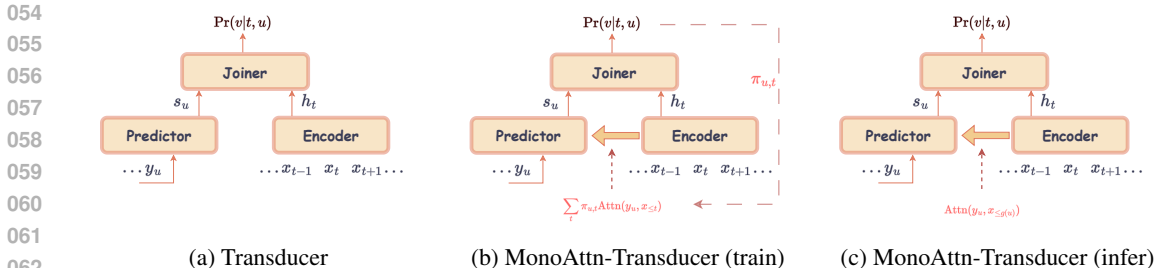


Figure 1: Illustration of MonoAttn-Transducer. During inference, the predictor state can attend to all generated encoder states through cross-attention. During training, the scope of cross-attention is adjusted based on the posterior alignment $\pi_{u,t}$ derived from the model’s prediction (Equation 5).

ability to attend to the input stream history during decoding, making it hard to manage reorderings. To address this issue, recent research (Liu et al., 2021; Tang et al., 2023) has started to explore the incorporation of cross-attention mechanism to enhance the capacity for handling complex non-monotonicity. Despite these efforts, the integration of cross-attention presents significant challenges. By integrating the predictor states with source history through attention, the representation of predictor states becomes relevant not only to the encoder states but also to the specific READ/WRITE path history (Tang et al., 2023). This results in an exponentially large state space for Transducer, hindering the application of the forward-backward algorithm for end-to-end training.

In this research, we present an efficient training algorithm for Transducer models to learn the monotonic cross-attention mechanism. This allows Transducer’s predictor to access source history in real-time inference (Figure 1c), improving its ability to handle tasks with non-monotonic alignments. As illustrated in Figure 1b, we leverage the forward-backward algorithm to infer the posterior probability of alignments between predictor and encoder states in training. This derived posterior alignment enables the estimation of context representation for each predictor state using expected soft attention. In this way, Transducer models adaptively adjust the scope of attention based on their predictions, avoiding the need to enumerate the exponentially large alignment space during training.

We conduct experiments on both speech-to-text and speech-to-speech simultaneous translation to demonstrate the generality of our approach across various modalities. MonoAttn-Transducer shows significant improvements in generation quality without a noticeable increase in latency in both *ideal* and *computation-aware* settings (§5). Further analysis reveals that MonoAttn-Transducer is particularly effective in handling samples with higher levels of non-monotonicity (§6).

2 BACKGROUND

2.1 STREAMING GENERATION

Streaming generation models typically process a streaming input $\mathbf{x} = \{x_1, \dots, x_T\}$ and generate a target sequence $\mathbf{y} = \{y_1, \dots, y_U\}$ in a streaming manner. To measure the amount of source information utilized during generation, a monotonic non-decreasing function $g(u)$ is introduced to represent the number of observed source tokens at the time of generating y_u .

2.2 TRANSDUCER

As shown in Figure 1a, Transducer model (Graves, 2012) comprises three components: an encoder, a predictor, and a joiner. The encoder unidirectionally encodes the received input prefix $x_{1:t}$ into a context representation h_t . The predictor functions similarly to an autoregressive language model, encoding the dependencies between tokens in the generated prefix $y_{1:u}$ into s_u . The joiner makes predictions based on the current source representation h_t and target representation s_u . If the model needs to READ more information to update the source representation for continued generation, a blank token ϵ is generated. Otherwise, a WRITE operation is performed, and the generated token is fed back into the predictor to obtain a new target representation. Each time h_t or s_u is updated, the joiner performs a prediction step until the entire source has been processed. The encoder and

108 predictor are usually modeled using either a recurrent neural network (Graves, 2012) or Transformer
 109 layers (Zhang et al., 2020). The joiner is typically composed of a feed-forward network.

110
 111 Since explicit alignment information for parallel pairs is not available during training, it is necessary
 112 to solve for the total probabilities of all READ/WRITE paths that can generate the target to perform
 113 maximum likelihood estimation. Given that the state space of Transducer form a two-dimensional
 114 lattice, the forward-backward algorithm can be utilized to compute the total probability. Define the
 115 forward and backward variables as:

$$\begin{aligned} \alpha(t, u) &:= \Pr(y_{1:u} | x_{1:t}) \\ \beta(t, u) &:= \Pr(y_{u+1:U} | x_{t:T}) \end{aligned} \quad (1)$$

118 The forward and backward variables for all $1 \leq t \leq T$ and $0 \leq u \leq U$ can be calculated recursively:

$$\begin{aligned} \alpha(t, u) &= \alpha(t-1, u) \Pr(\epsilon | t-1, u) + \alpha(t, u-1) \Pr(y_u | t, u-1) \\ \beta(t, u) &= \beta(t+1, u) \Pr(\epsilon | t, u) + \beta(t, u+1) \Pr(y_{u+1} | t, u) \end{aligned} \quad (2)$$

122 with initial condition $\alpha(1, 0) = 1$ and $\beta(T, U) = \Pr(\epsilon | T, U)$. $\Pr(v | t, u)$ denotes the probability of
 123 generating token v from h_t and s_u , $v \in \mathcal{V} \cup \{\epsilon\}$. The total output probability is:

$$\Pr(\mathbf{y} | \mathbf{x}) = \alpha(T, U) \Pr(\epsilon | T, U). \quad (3)$$

126 By leveraging the forward-backward algorithm, Transducer models are trained to implicitly acquire
 127 the READ/WRITE policy from the data.

129 3 METHOD

131 In this section, we provide a detailed introduction to our proposed MonoAttn-Transducer.

133 3.1 OVERVIEW

135 MonoAttn-Transducer works similarly to standard Transducer, with the key difference being that its
 136 predictor can attend to the encoder history using monotonic attention. During streaming generation,
 137 the scope of monotonic attention includes all source context representations that have already ap-
 138 peared. Formally, when the predictor encodes the u -th target state, it depends on the representations
 139 of previous target states and the existing source context:

$$s_u = f_\theta(s_{0:u-1}, h_{1:g(u)}), \quad (4)$$

142 where $1 \leq u \leq U$ and $g(u)$ denotes the number of observed source tokens at the time of generating
 143 y_u . The edge case s_0 is defined as $s_0 = f_\theta(h_1)$. In both Transducer and MonoAttn-Transducer,
 144 token y_u is generated based on source representation $h_{g(u)}$ and target representation s_{u-1} . Given s_{u-1}
 145 can only attend to source contexts up to $g(u-1)$ through monotonic attention, related information
 146 in $x_{g(u-1)+1:g(u)}$ should ideally be encoded within $h_{g(u)}$.

148 3.2 TRAINING ALGORITHM

149 Training MonoAttn-Transducer is challenging as it exponentially expands Transducer’s state space.
 150 To address this issue, we firstly leverage the forward-backward algorithm to compute the posterior
 151 probability of aligning target representation s_u with source representation h_t (i.e., the probability of
 152 generating token y_u immediately after reading x_t). This posterior alignment is then used to estimate
 153 the expected context vector in the monotonic cross-attention for each predictor state in training.
 154 Detailed explanations are provided in the following.

156 3.2.1 POSTERIOR ALIGNMENT

157 Suppose we have a probability lattice $\Pr(v | t, u)$, representing the probability of generating token
 158 v from h_t and s_u , for $1 \leq t \leq T$, $0 \leq u \leq U$, and $v \in \mathcal{V} \cup \{\epsilon\}$. The posterior probability of
 159 generating y_u at the moment x_t is read can be represented by:

$$\pi_{u,t} = \frac{\Pr(y_{1:u-1} | x_{1:t}) \Pr(y_u | t, u-1) \Pr(y_{u+1:U} | x_{t:T})}{\Pr(y_{1:U} | x_{1:T})} \quad (5)$$

with the edge case:

$$\pi_{0,t} = \begin{cases} 1 & t = 1 \\ 0 & t \neq 1 \end{cases} \quad (6)$$

which implies that the predictor state s_0 is generated immediately after the first source token arrives. Using the forward and backward variables introduced in Section 2.2, Equation 5 can be concisely expressed as follows:

$$\pi_{u,t} = \frac{\alpha(t, u - 1)\Pr(y_u|t, u - 1)\beta(t, u)}{\alpha(T, U)\Pr(\epsilon|T, U)}. \quad (7)$$

This guarantees that the posterior alignment probability for all pairs (t, u) can be solved in $O(TU)$ time using the above forward-backward algorithm, facilitating the calculation of the expected context representation introduced later.

3.2.2 MONOTONIC ATTENTION

The incorporation of monotonic attention makes the representation of predictor states relevant to specific READ/WRITE history, leading to a prohibitively large state space for enumerating alignments. Therefore, we turn to estimate the context vector in monotonic attention based on the posterior alignment probability during training. This approach enables the model to adaptively adjust the scope of cross-attention according to its prediction. Consequently, MonoAttn-Transducer learns a monotonic attention mechanism while maintaining the same time and space complexity as Transducer.

Formally, given the energy $e_{u,t}$ for the pair consisting of encoder state h_t and predictor state s_u , as well as the posterior alignment probability $\pi_{u,t}$, the expected context representation c_u for predictor state s_u can be expressed as:

$$c_u = \sum_{t=1}^T \pi_{u,t} \sum_{t'=1}^t \frac{\exp(e_{u,t'})}{\sum_{t''=1}^t \exp(e_{u,t''})} h_{t'}. \quad (8)$$

This indicates that the expected context representation c_u is a weighted sum of context representations under various amount of source information, with the weights given by the posterior alignment probability $\pi_{u,t}$. The nested summation operations in Equation 8 may lead to an increase in computational complexity. Fortunately, Arivazhagan et al. (2019) suggests that it can be rewritten as:

$$\begin{aligned} \phi_{u,t} &= \sum_{t'=t}^T \frac{\pi_{u,t'} \exp(e_{u,t})}{\sum_{t''=1}^{t'} \exp(e_{u,t''})} \\ c_u &= \sum_{t=1}^T \phi_{u,t} h_t \end{aligned} \quad (9)$$

Equation 9 can then be computed efficiently using cumulative sum operations (Arivazhagan et al., 2019).

3.2.3 TRAINING WITH PRIOR ALIGNMENT

The above algorithm facilitates MonoAttn-Transducer in learning monotonic cross-attention with posterior alignment probability. However, this presents a chicken-and-egg paradox: the posterior alignment is derived from an output probability lattice constructed using an estimated context representation, while the context vector is, in turn, estimated using a posterior alignment. We address this problem by using a prior alignment to construct a prior output probability lattice. This lattice is then used to infer the posterior alignment and train MonoAttn-Transducer’s monotonic attention.

There are several options for the prior alignment $p_{u,t}$. The simplest one is the uniform distribution, which assigns an equal probability of being generated at any timestep for all the target tokens:

$$p_{u,t}^{\text{uni}} = \frac{1}{T}, \quad 1 \leq t \leq T, \quad 1 \leq u \leq U. \quad (10)$$

The edge case $p_{0,t}^{\text{uni}}$ is similar to the situation of $\pi_{0,t}$, where all the probability mass is concentrated at $t = 1$.

Algorithm 1 Training Algorithm of MonoAttn-Transducer**Input:** Source x , Target y , Chunk Size C **Output:** Training Loss \mathcal{L}

- 1: Compute prior alignment $p_{u,t}^{\text{dia}}$ (Eq. 11)
- 2: Compute chunk-synchronized prior alignment $\tilde{p}_{u,t}^{\text{dia}}$ based on chunk size C (Eq. 12)
- 3: Estimate context c_u^{prior} with $\tilde{p}_{u,t}^{\text{dia}}$ (Eq. 9)
- 4: Forward MonoAttn-Transducer with c_u^{prior}
- 5: Infer posterior alignment $\pi_{u,t}$ (Eq. 7)
- 6: Compute chunk-synchronized posterior alignment $\tilde{\pi}_{u,t}$ based on chunk size C (Eq. 12)
- 7: Estimate context c_u with $\tilde{\pi}_{u,t}$ (Eq. 9)
- 8: Forward MonoAttn-Transducer with c_u
- 9: Calculate total output probability \mathcal{L} (Eq. 3)
- 10: **return** \mathcal{L}

However, it is preferable to select a more reasonable prior. An ideal prior alignment should ensure that the posterior alignment, derived from the lattice constructed using the prior, can accurately estimate the expected context representation. In streaming generation tasks, even though there may be reorderings in the mapping from source to target, a certain level of monotonic alignment is generally maintained. Therefore, we propose introducing a prior distribution $p_{u,t}^{\text{dia}}$, which assumes that the number of tokens generated for each READ action is uniformly distributed:

$$w_{u,t} = \exp\left(-\left|u - \frac{t \cdot U}{T}\right|\right)$$

$$p_{u,t}^{\text{dia}} = \frac{w_{u,t}}{\sum_{t'=1}^T w_{u,t'}} \quad (11)$$

for $1 \leq t \leq T$, $1 \leq u \leq U$. The edge case $p_{0,t}^{\text{dia}}$ is handled in the same manner as $p_{0,t}^{\text{uni}}$. This prior assumes a uniform mapping between the source and target, such that each target token is most likely generated at the time its corresponding source token is read. The probability decreases as the time difference from this moment increases.

In the following, we will use $p_{0,t}^{\text{dia}}$ as the default choice for prior alignment and compare the differences between using $p_{0,t}^{\text{uni}}$ and $p_{0,t}^{\text{dia}}$ in the ablation study (Section 6.1).

3.2.4 CHUNK SYNCHRONIZATION

In speech audio, there often exists strong temporal dependencies between adjacent frames. Therefore, a chunk size C is typically set, and the streaming model makes decisions only after receiving a speech chunk (Ma et al., 2020c). In terms of Transducer models, when a READ action is taken, the source representation is updated after a new speech chunk is read. The new source representation is then set as the representation of the last frame in the chunk (Liu et al., 2021; Tang et al., 2023). In such a situation, the receptive field of MonoAttn-Transducer’s cross-attention for predictor state s_u encompasses all hidden states in the received chunks, i.e., $h_{1:C \cdot \tilde{g}(u)}$, where $\tilde{g}(u)$ denotes the number of received chunks when generating token y_u . To bridge the gap between training and inference, the posterior alignment probability utilized in training process is adjusted by transferring all the probability mass on encoder states within a chunk to the last state in the chunk:

$$\tilde{\pi}_{u,t} = \begin{cases} \sum_{t'=(d-1) \cdot C+1}^{d \cdot C} \pi_{u,t'} & t = d \cdot C \\ 0 & t \neq d \cdot C \end{cases} \quad (12)$$

for $d = 1, 2, 3, \dots$

The prior alignment probability is adjusted in the same manner. We detail the entire training process in Algorithm 1.

4 RELATED WORK

Our work is closely related to researches in designing cross-attention modules for Transducer models. Prabhavalkar et al. (2017) pioneered the use of attention to link the predictor and encoder.

Table 1: Comparison of Transducer-based streaming models. *Computational Complexity* refers to the number of forward passes executed by the predictor in inference. *Memory Overhead* refers to the memory consumption of the attention module in training.

Method	Merge Module	Computational Complexity	Memory Overhead
Transducer (Graves, 2012)	Joiner	$O(U)$	N/A
CAAT (Liu et al., 2021)	Joiner	$O(U)$	$O(T)$
TAED (Tang et al., 2023)	Predictor, Joiner	$O(U + T)$	$O(T)$
MonoAttn-Transducer (Ours)	Predictor, Joiner	$O(U)$	$O(1)$

However, their design requires the entire source to be available, limiting it to offline generation. For streaming generation, the receptive field of attention must synchronize with the input. This synchronization leads to an exponentially large state space, which significantly complicates the training process. To mitigate this issue, Liu et al. (2021) separated the predictor’s cross-attention from its self-attention, ensuring that cross-attention occurs only after self-attention. This approach maintains the independence of predictor states from READ/WRITE path history, allowing for standard training methods. However, this separation limits the richness of the predictor’s learned representations. Alternatively, Tang et al. (2023) proposed updating the representation of all predictor states whenever a new source token is received. While this method also preserves the independence of predictor states from READ/WRITE path history, it significantly increases both inference-time computational complexity and training-time memory requirements. It necessitates an additional $(T - 1)$ forward passes of the predictor during decoding, which adversely affects latency-sensitive streaming generation. Furthermore, the GPU memory usage for attention during training increases from $O(1)$ to $O(T)$, leading to prohibitively high training costs and limiting the model’s scalability. In contrast to the above, the proposed MonoAttn-Transducer maintains the same time complexity and memory overhead as Transducer. A detailed comparison between these methods is summarized in Table 1.

Our work is also related to researches in designing attention modules for streaming AED models. These works often introduce Bernoulli variables to indicate READ/WRITE actions. The distribution of these variables is used to estimate monotonic alignment and to compute the expected context representation in training (Raffel et al., 2017). Depending on the setting of attention window, these works can be classified into monotonic hard attention (Raffel et al., 2017), monotonic chunkwise attention (MoChA; Chiu & Raffel, 2018), and monotonic infinite lookback attention (MILk; Arivazhagan et al., 2019). Ma et al. (2020d) subsequently introduced the MILk mechanism to Transformer models, and Ma et al. (2023b) further proposed a numerically-stable algorithm for estimating monotonic alignment. Unlike the aforementioned works, our approach learns monotonic attention based on the posterior alignment of Transducer, avoiding the use of unstable Bernoulli variables.

5 EXPERIMENTS

We validate the performance of our MonoAttn-Transducer on two typical streaming generation tasks: speech-to-text and speech-to-speech simultaneous translation. The differences in grammatical structures between the source and target languages often necessitate word reordering during generating translation. This property makes the simultaneous translation task well-suited for evaluating the ability of MonoAttn-Transducer in handling non-monotonic alignments.

5.1 EXPERIMENTAL SETUP

Datasets We conduct experiments on two language pairs of MuST-C speech-to-text translation datasets: English to German (En→De) and English to Spanish (En→Es) (Di Gangi et al., 2019). For speech-to-speech experiments, we evaluate models on CVSS-C French to English (Fr→En) dataset (Jia et al., 2022).

Model Configuration We use the open-source implementation of Transformer-Transducer (Zhang et al., 2020) from Liu et al. (2021) as baseline and build our MonoAttn-Transducer upon it. The speech encoder consists of two layers of causal 2D-convolution followed by 16 chunk-wise Trans-

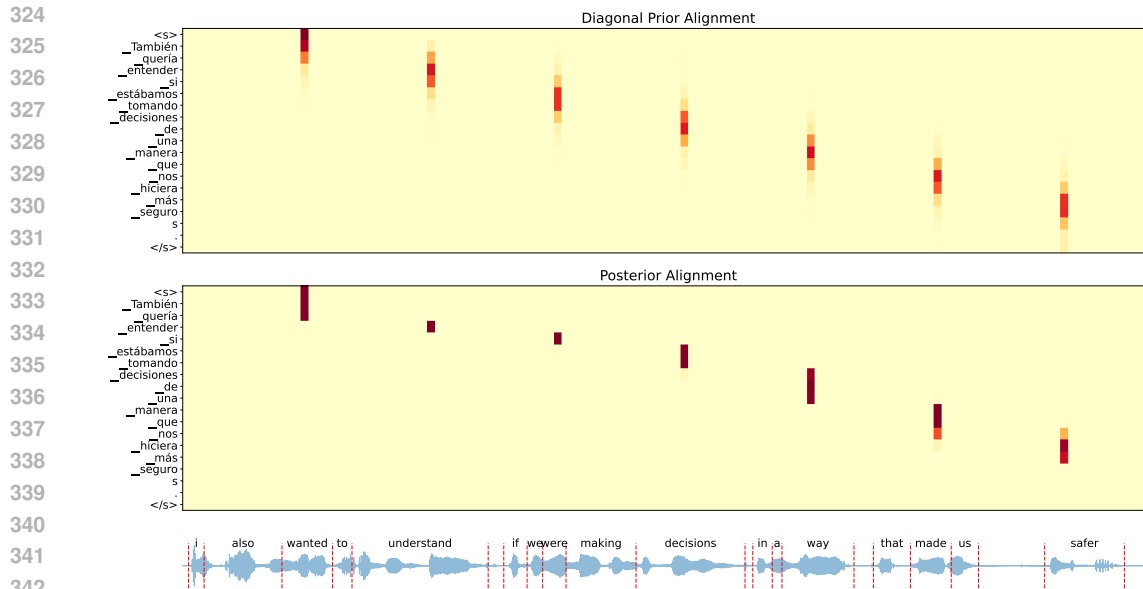


Figure 2: An example of diagonal prior and posterior alignment from MuST-C English-to-Spanish training corpus. The vertical axis represents the target subword sequence and the horizontal axis represents the speech waveform. Darker areas indicate higher alignment probabilities. Chunk size in this example is set to 640ms. More examples are provided in Appendix D.

former layers with pre-norm. Each convolution layer has a 3×3 kernel with 64 channels and a stride size of 2, resulting in a downsampling ratio of 4. In chunk-wise Transformer layers, the speech encoder can access states from all previous chunks and one chunk ahead of the current chunk (Wu et al., 2020; Shi et al., 2021). The chunk size is adjusted within the set $\{320, 640, 960, 1280\}$ ms. Offline results are obtained by setting the chunk size longer than any utterance in the corpus. Both sinusoidal positional encoding (Vaswani et al., 2017) and relative positional attention (Shaw et al., 2018) are incorporated into the speech encoder. Sinusoidal positional encoding is applied after the convolution layers. The predictor comprises two autoregressive Transformer layers with post-norm, utilizing only sinusoidal positional encoding. The monotonic attention is similar to standard cross-attention but differs in its receptive field. The joiner is implemented as a simple FFN. We incorporate the multi-step decision mechanism (Liu et al., 2021) with a decision step of 4. All Transformer layers described above are configured with a 512 embedding dimension, 8 attention heads and a 2048 FFN dimension. The total number of parameters for the Transducer baseline and MonoAttn-Transducer are 65M and 67M, respectively. More implementation details are provided in Appendix A.

Evaluation We use SimulEval toolkit (Ma et al., 2020b) for evaluation. Translation quality is assessed using case-sensitive detokenized BLEU (Papineni et al., 2002; Post, 2018) and neural-based COMET-22 score. Latency is measured by word-level Average Lagging (AL; Ma et al., 2019; 2020c).¹ For speech-to-speech experiments, translation quality is assessed using ASR-BLEU and latency is measured by delay of generated waveform chunks (Ma et al., 2022).

5.2 MAIN RESULTS

We evaluate the performance of MonoAttn-Transducer against Transducer baseline across various latency conditions obtained by varying the chunk size. In this comparison, we consider two configurations of MonoAttn-Transducer. The first, referred to as MonoAttn-Transducer-*Posterior*, is trained strictly according to Algorithm 1. The second, termed MonoAttn-Transducer-*Prior*, is optimized

¹Numerical results with more metrics are provided in Appendix C. Notably, Table 7 presents a comparison of the *computation-aware* latency metrics for AL and LAAL (Papi et al., 2022) between the Transducer and MonoAttn-Transducer models.

Table 2: Comparison of MonoAttn-Transducer and Transducer across various chunk size settings on MuST-C English to German and English to Spanish datasets.

		<i>En-Es</i>					<i>En-De</i>				
		320	640	960	1280	∞	320	640	960	1280	∞
Transducer	AL (<i>ms</i> , ↓)	886	1193	1591	1997	-	1126	1434	1830	2215	-
	BLEU (↑)	24.33	25.82	26.36	26.40	26.75	19.99	22.10	22.20	22.96	23.10
	COMET (↑)	67.94	69.92	70.48	70.65	71.14	62.81	65.01	65.75	66.26	67.03
MonoAttn-Transducer (<i>Posterior</i>)	AL (<i>ms</i> , ↓)	997	1239	1606	1991	-	1215	1470	1860	2215	-
	BLEU (↑)	24.72	26.74	27.05	27.41	27.48	20.22	22.47	22.94	23.74	24.42
	COMET (↑)	68.98	70.71	71.21	71.90	72.24	64.24	67.06	68.22	68.54	69.82
MonoAttn-Transducer (<i>Prior</i>)	AL (<i>ms</i> , ↓)	932	1182	1599	1967	-	1138	1413	1826	2191	-
	BLEU (↑)	23.00	26.46	27.07	27.42	27.48	19.26	22.62	23.51	24.01	24.42
	COMET (↑)	68.24	70.45	71.33	71.99	72.24	63.85	67.63	68.65	69.27	69.82

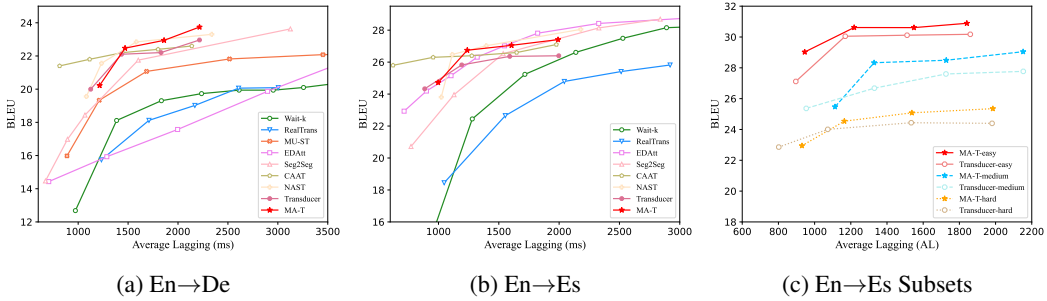


Figure 3: (a), (b): Results of translation quality (BLEU) against latency (Average Lagging, AL) on MuST-C English to German and English to Spanish datasets. (c): Performance on MuST-C English to Spanish test subsets categorized by non-monotonicity. In the figures above, *MA-T* denotes MonoAttn-Transducer.

directly using prior alignment, without inferring the posterior (calculate total output probability \mathcal{L} using c_u^{prior}). Results are shown in Table 2.

It can be observed that MonoAttn-Transducer-*Posterior* significantly outperforms the Transducer baseline across various settings of chunk size in both translation directions. Specifically, in En-Es, it shows an average improvement of 0.75 BLEU or 0.95 COMET score in generation quality under different latency conditions. In En-De, it achieves an even more significant improvement, with an average increase of as much as 2.06 COMET score, while latency remains nearly unchanged. Further analysis reveals that the benefits of learning monotonic attention are more pronounced with a larger chunk size. Notably, in scenarios where latency exceeds 1.5s and during offline generation, the average improvement reaches 0.88 BLEU or 1.77 COMET score. This can be attributed to MonoAttn-Transducer benefiting more from monotonic attention to handle reorderings when it has flexibility to wait for more source information.

Moreover, we have observed some notable results of MonoAttn-Transducer-*Prior*. With a larger chunk size, the performance of MonoAttn-Transducer-*Prior* is comparable to that of MonoAttn-Transducer-*Posterior*, and even slightly outperforming the latter in En-De. However, there exists a significant performance drop with a smaller chunk size. Specifically, with a chunk size of 320ms, MonoAttn-Transducer-*Prior*'s generation quality is on average 1.03 BLEU lower than Transducer baseline under similar latency conditions. This phenomenon highlights the importance of learning monotonic attention through inferring posterior alignment. From the chunk synchronization mechanism described in Equation 12, smaller chunk sizes require finer alignment granularity between the predictor and encoder states. This increased granularity necessitates more precise alignment to estimate the expected context representation during training. Figure 2 provides an example of diagonal prior and posterior alignment. While the diagonal prior generally captures the trend of the alignment information, it can be skewed by the uneven distribution of speech information and possible local reorderings. In contrast, the inferred posterior offers a more confident and accurate alignment prob-

Table 3: Performance on CVSS-C French to English speech-to-speech translation.

	Chunk Size (ms)	320	Offline
Ours	ASR-BLEU (\uparrow)	18.3	19.3
	AL (ms, \downarrow)	118	-
Transducer	ASR-BLEU (\uparrow)	17.1	18.0
	AL (ms, \downarrow)	153	-

Table 4: Performance of MonoAttn-Transducer with different choices of prior alignment.

	Chunk Size (ms)	320	640	960	1280
p^{dia}	BLEU (\uparrow)	24.72	26.74	27.05	27.41
	AL (ms, \downarrow)	997	1239	1606	1991
p^{uni}	BLEU (\uparrow)	24.89	26.68	27.26	27.11
	AL (ms, \downarrow)	993	1249	1601	1983

ability. For instance, the diagonal prior assigns a high probability to aligning the word “*si (if)*” with the timestep preceding the waveform of “*if*”, while the inferred posterior corrects this misalignment. Therefore, learning monotonic attention with posterior alignment leads to a more accurate estimation of context representation and improved performance.² In subsequent experiments, we represent MonoAttn-Transducer using the results of MonoAttn-Transducer-*Posterior*.

5.3 COMPARISON WITH STATE-OF-THE-ART

We compare MonoAttn-Transducer with state-of-the-art open-source approaches in simultaneous translation, including Wait- k (Ma et al., 2020c), RealTrans (Zeng et al., 2021), CAAT (Liu et al., 2021), MU-ST (Zhang et al., 2022), EDAtt (Papi et al., 2023), Seg2Seg (Zhang & Feng, 2023) and NAST (Ma et al., 2024). Further details about baselines are available in Appendix B. Results are plotted in Figure 3a and 3b. We observe that learning monotonic attention significantly enhances the performance of Transducer, making it comparable to state-of-the-art models. Compared to CAAT, another Transducer-based model, MonoAttn-Transducer demonstrates superiority in scenarios with less stringent latency requirements. Under a latency of approximately 2s, it outperforms CAAT by 1.1 BLEU in En-De. This clearly demonstrates the advantage of MonoAttn-Transducer’s tightly coupled self-attention and cross-attention modules in the predictor, which facilitates the learning of richer representations.

As discussed in Section 4, TAED is another Transducer-based model highly relevant to our work. However, the code and distilled data used to train TAED in Tang et al. (2023) have not been made publicly available. This lack of open access hinders a fair comparison of TAED with our MonoAttn-Transducer. Despite this, we attempt to analyze the performance by comparing each with Transducer baseline in their respective experimental settings. The comparison is shown in Table 8. We have observed that the improvement from TAED is more pronounced with smaller chunk sizes, which contrasts with the results of MonoAttn-Transducer. We speculate that this is because, in TAED, the representations of all generated predictor states are updated every time the encoder receives a new speech chunk. This helps TAED generate more accurate representations when the chunk size is small. However, this mechanism in TAED incurs an $O(T + U)$ forward propagation cost during simultaneous inference, which can significantly increase latency in practice due to heavy computational overhead when the chunk size is small. **In contrast, MonoAttn-Transducer maintains an $O(U)$ complexity as Transducer baseline. As shown in Table 7, this property minimizes the gap between ideal and computation-aware latency, offering advantages in real-time applications.**

5.4 RESULTS OF SPEECH GENERATION

Speech-to-speech simultaneous translation requires implicitly performing ASR, MT and TTS simultaneously, and also handling the non-monotonic alignments between languages, making it suitable to evaluate models on streaming speech generation. We adopted a *textless* setup in our experiments, directly modeling the mapping between speech (Zhao et al., 2024). Results are provided in Table 3.

The results demonstrate that MonoAttn-Transducer significantly reduces generation latency (AL). With a chunk size of 320ms, it achieves Transducer’s offline generation quality, but reducing lagging to 118ms. For offline settings, our approaches further improves speech generation quality (19.3 vs. 18.0). These results highlight the effectiveness of our approach in achieving a better quality-latency trade-off also for streaming speech generation.

²We present a comparison between the prior and posterior under various chunk sizes in Appendix D. A key observation is that as the alignment granularity becomes finer, the differences gradually increases.

486 6 ANALYSIS

487 6.1 CHOICE OF PRIOR ALIGNMENT

488 In Section 3.2.3, we introduced two choices for prior alignment: the uniform prior p^{uni} , which
 489 assumes an equal probability of generation at each time step; and the diagonal prior p^{dia} , which
 490 prefers ideal synchrony between the source and target. We employed the diagonal prior p^{dia} as the
 491 default choice in the aforementioned experiments. In this section, we examine the impact of differ-
 492 ent choices. The results are displayed in Table 4. As shown, MonoAttn-Transducer’s performance
 493 demonstrates robustness to the choice of prior alignment, with only minor impacts on both transla-
 494 tion quality and latency across all chunk size settings. In Appendix D, we visualize the posterior
 495 alignment when using *different* priors. We have observed that, even with significant differences in
 496 the prior distribution, the posterior remains fairly robust when the chunk size is constant. This nice
 497 property reinforces the robustness of using the inferred posterior to train monotonic attention.
 498
 499

500 6.2 HANDLING NON-MONOTONICITY

501 To illustrate MonoAttn-Transducer’s capability in handling reorderings through learning monotonic
 502 attention, we evaluate its performance against the Transducer baseline across samples with varying
 503 levels of non-monotonicity. Intuitively, samples with a higher number of crosses in the alignments
 504 between source transcription and reference text pose greater challenges. We therefore evenly par-
 505 tition the test set based on the number of cross-alignments, labeling them as easy, medium and
 506 hard.³ The results are presented in Figure 3c. We observe that MonoAttn-Transducer shows a
 507 more substantial improvement over Transducer in the medium and hard subsets across most chunk
 508 size settings. However, with a chunk size of 320ms, the improvement is particularly notable in the
 509 easy subset. These findings highlight the unique capabilities of MonoAttn-Transducer in managing
 510 non-monotonic alignments. As analyzed in Section 5.2, MonoAttn-Transducer benefits more from
 511 learning monotonic attention with a larger chunk size, and this enhanced ability is evident in subsets
 512 with higher levels of non-monotonicity. On the other hand, when the chunk size is extremely small,
 513 MonoAttn-Transducer has limited flexibility to wait for more source information before processing,
 514 thus showing more significant improvement in the easy subset under the condition.
 515

516 6.3 TRAINING EFFICIENCY

517 **Training Time:** We analyze each step in Algorithm 1 to compare the training time differences
 518 between MonoAttn-Transducer and baseline. We observe that Lines 1, 2, 6 involve naive matrix
 519 computation without requiring gradients. The additional time overhead introduced by our method
 520 arises from Lines 3, 4, 5. Specifically, this includes an additional forward pass of the predictor and
 521 the computation for the posterior alignment. The overhead from the posterior calculation is approx-
 522 imately equivalent to that incurred during loss calculation, as both rely on the forward-backward
 523 algorithm. Empirically, we found MonoAttn-Transducer is 1.33 times slower than Transducer base-
 524 line with the same configuration on Nvidia L40 GPU.
 525

526 **Memory Consumption:** Compared to baseline, the additional memory overhead of MonoAttn-
 527 Transducer comes solely from its monotonic attention module. The extra forward pass of the predic-
 528 tor is performed without requiring gradients, so it is excluded from the computation graph. Empir-
 529 ically, we observed that the peak memory usage of Transducer baseline is 28GB, while MonoAttn-
 530 Transducer exhibits a slightly higher peak usage of 32GB when the total number of source frames
 531 is fixed at 40,000 on a single Nvidia L40 GPU.

532 7 CONCLUSION

533 In this paper, we propose an efficient algorithm for Transducer models to learn monotonic atten-
 534 tion. Extensive experiments demonstrate that our MonoAttn-Transducer significantly improves the
 535 ability in handling non-monotonic alignments in streaming generation, offering a robust solution for
 536 Transducer-based frameworks to tackle more complex streaming generation tasks.
 537
 538

539 ³The easy subset includes samples with a cross count of 1 or fewer. The medium subset contains samples
 with a cross count between 2 and 6. Samples with a cross count greater than 6 are classified as hard.

REFERENCES

- 540
541
542 Naveen Arivazhagan, Colin Cherry, Wolfgang Macherey, Chung-Cheng Chiu, Semih Yavuz, Ruoming Pang, Wei Li, and Colin Raffel. Monotonic infinite lookback attention for simultaneous machine translation. In Anna Korhonen, David Traum, and Lluís Màrquez (eds.), *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 1313–1323, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1126. URL <https://aclanthology.org/P19-1126>.
543
544
545
546
547
- 548 Dzmitry Bahdanau. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
549
550
- 551 Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. In Yoshua Bengio and Yann LeCun (eds.), *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. URL <http://arxiv.org/abs/1409.0473>.
552
553
554
- 555 Chung-Cheng Chiu and Colin Raffel. Monotonic chunkwise attention. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=Hko85p1CW>.
556
557
- 558 Kyunghyun Cho and Masha Esipova. Can neural machine translation do simultaneous translation? *CoRR*, abs/1606.02012, 2016. URL <http://arxiv.org/abs/1606.02012>.
559
560
- 561 Shun-Po Chuang, Yung-Sung Chuang, Chih-Chiang Chang, and Hung-yi Lee. Investigating the re-ordering capability in CTC-based non-autoregressive end-to-end speech translation. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pp. 1068–1077, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.findings-acl.92. URL <https://aclanthology.org/2021.findings-acl.92>.
562
563
564
565
566
- 567 Mattia A. Di Gangi, Roldano Cattoni, Luisa Bentivogli, Matteo Negri, and Marco Turchi. MuST-C: a Multilingual Speech Translation Corpus. In Jill Burstein, Christy Doran, and Tamar Solorio (eds.), *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)*, pp. 2012–2017, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1202. URL <https://aclanthology.org/N19-1202>.
568
569
570
571
572
- 573 Alex Graves. Sequence transduction with recurrent neural networks. *CoRR*, abs/1211.3711, 2012. URL <http://arxiv.org/abs/1211.3711>.
574
575
- 576 Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber. Connectionist temporal classification: Labelling unsegmented sequence data with recurrent neural networks. In *Proceedings of the 23rd International Conference on Machine Learning, ICML '06*, pp. 369–376, New York, NY, USA, 2006. Association for Computing Machinery. ISBN 1595933832. doi: 10.1145/1143844.1143891. URL <https://doi.org/10.1145/1143844.1143891>.
577
578
579
580
- 581 Jiatao Gu, Graham Neubig, Kyunghyun Cho, and Victor O.K. Li. Learning to translate in real-time with neural machine translation. In Mirella Lapata, Phil Blunsom, and Alexander Koller (eds.), *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pp. 1053–1062, Valencia, Spain, April 2017. Association for Computational Linguistics. URL <https://aclanthology.org/E17-1099>.
582
583
584
585
- 586 Ye Jia, Michelle Tadmor Ramanovich, Quan Wang, and Heiga Zen. CVSS corpus and massively multilingual speech-to-speech translation. In *Proceedings of Language Resources and Evaluation Conference (LREC)*, pp. 6691–6703, 2022.
587
588
589
- 590 Yoon Kim and Alexander M. Rush. Sequence-level knowledge distillation. In Jian Su, Kevin Duh, and Xavier Carreras (eds.), *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp. 1317–1327, Austin, Texas, November 2016. Association for Computational Linguistics. doi: 10.18653/v1/D16-1139. URL <https://aclanthology.org/D16-1139>.
591
592
593

- 594 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua
595 Bengio and Yann LeCun (eds.), *3rd International Conference on Learning Representations, ICLR*
596 *2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. URL <http://arxiv.org/abs/1412.6980>.
- 597
598
- 599 Taku Kudo and John Richardson. SentencePiece: A simple and language independent subword
600 tokenizer and detokenizer for neural text processing. In Eduardo Blanco and Wei Lu (eds.),
601 *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing:*
602 *System Demonstrations*, pp. 66–71, Brussels, Belgium, November 2018. Association for Com-
603 putational Linguistics. doi: 10.18653/v1/D18-2012. URL <https://aclanthology.org/D18-2012>.
- 604
- 605 Dan Liu, Mengge Du, Xiaoxi Li, Ya Li, and Enhong Chen. Cross attention augmented transducer
606 networks for simultaneous translation. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia,
607 and Scott Wen-tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Nat-
608 ural Language Processing*, pp. 39–55, Online and Punta Cana, Dominican Republic, November
609 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.4. URL
610 <https://aclanthology.org/2021.emnlp-main.4>.
- 611
- 612 Mingbo Ma, Liang Huang, Hao Xiong, Renjie Zheng, Kaibo Liu, Baigong Zheng, Chuanqiang
613 Zhang, Zhongjun He, Hairong Liu, Xing Li, Hua Wu, and Haifeng Wang. STACL: Simultaneous
614 translation with implicit anticipation and controllable latency using prefix-to-prefix framework.
615 In Anna Korhonen, David Traum, and Lluís Màrquez (eds.), *Proceedings of the 57th Annual*
616 *Meeting of the Association for Computational Linguistics*, pp. 3025–3036, Florence, Italy, July
617 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1289. URL <https://aclanthology.org/P19-1289>.
- 618
- 619 Mingbo Ma, Baigong Zheng, Kaibo Liu, Renjie Zheng, Hairong Liu, Kainan Peng, Kenneth
620 Church, and Liang Huang. Incremental text-to-speech synthesis with prefix-to-prefix frame-
621 work. In Trevor Cohn, Yulan He, and Yang Liu (eds.), *Findings of the Association for Com-
622 putational Linguistics: EMNLP 2020*, pp. 3886–3896, Online, November 2020a. Association
623 for Computational Linguistics. doi: 10.18653/v1/2020.findings-emnlp.346. URL <https://aclanthology.org/2020.findings-emnlp.346>.
- 624
- 625 Xutai Ma, Mohammad Javad Dousti, Changhan Wang, Jiatao Gu, and Juan Pino. SIMULEVAL:
626 An evaluation toolkit for simultaneous translation. In Qun Liu and David Schlangen (eds.), *Pro-
627 ceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System*
628 *Demonstrations*, pp. 144–150, Online, October 2020b. Association for Computational Linguis-
629 tics. doi: 10.18653/v1/2020.emnlp-demos.19. URL [https://aclanthology.org/2020.](https://aclanthology.org/2020.emnlp-demos.19)
630 [emnlp-demos.19](https://aclanthology.org/2020.emnlp-demos.19).
- 631
- 632 Xutai Ma, Juan Pino, and Philipp Koehn. SimulMT to SimulST: Adapting simultaneous text trans-
633 lation to end-to-end simultaneous speech translation. In Kam-Fai Wong, Kevin Knight, and Hua
634 Wu (eds.), *Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for*
635 *Computational Linguistics and the 10th International Joint Conference on Natural Language*
636 *Processing*, pp. 582–587, Suzhou, China, December 2020c. Association for Computational Lin-
637 guistics. URL <https://aclanthology.org/2020.aacl-main.58>.
- 638
- 639 Xutai Ma, Juan Miguel Pino, James Cross, Liezl Puzon, and Jiatao Gu. Monotonic multihead
640 attention. In *International Conference on Learning Representations*, 2020d. URL [https://](https://openreview.net/forum?id=Hyg96gBKPS)
641 openreview.net/forum?id=Hyg96gBKPS.
- 642
- 643 Xutai Ma, Hongyu Gong, Danni Liu, Ann Lee, Yun Tang, Peng-Jen Chen, Wei-Ning Hsu, Phillip
644 Koehn, and Juan Pino. Direct simultaneous speech-to-speech translation with variational mono-
645 tonic multihead attention, 2022.
- 646
- 647 Xutai Ma, Anna Sun, Siqu Ouyang, Hirofumi Inaguma, and Paden Tomasello. Efficient monotonic
multihead attention. *arXiv preprint arXiv:2312.04515*, 2023a.
- Xutai Ma, Anna Sun, Siqu Ouyang, Hirofumi Inaguma, and Paden Tomasello. Efficient monotonic
multihead attention, 2023b. URL <https://arxiv.org/abs/2312.04515>.

- 648 Zhengrui Ma, Shaolei Zhang, Shoutao Guo, Chenze Shao, Min Zhang, and Yang Feng. Non-
649 autoregressive streaming transformer for simultaneous translation. In Houda Bouamor, Juan Pino,
650 and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Lan-
651 guage Processing*, pp. 5177–5190, Singapore, December 2023c. Association for Computational
652 Linguistics. doi: 10.18653/v1/2023.emnlp-main.314. URL [https://aclanthology.org/
653 2023.emnlp-main.314](https://aclanthology.org/2023.emnlp-main.314).
- 654 Zhengrui Ma, Qingkai Fang, Shaolei Zhang, Shoutao Guo, Yang Feng, and Min Zhang. A
655 non-autoregressive generation framework for end-to-end simultaneous speech-to-any transla-
656 tion. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd An-
657 nual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp.
658 1557–1575, Bangkok, Thailand, August 2024. Association for Computational Linguistics. URL
659 <https://aclanthology.org/2024.acl-long.85>.
- 660 Michael McAuliffe, Michaela Socolof, Sarah Mihuc, Michael Wagner, and Morgan Sonderegger.
661 Montreal forced aligner: Trainable text-speech alignment using kaldi. In *Interspeech 2017*, pp.
662 498–502, 2017. doi: 10.21437/Interspeech.2017-1386.
- 663 Sara Papi, Marco Gaido, Matteo Negri, and Marco Turchi. Over-generation cannot be rewarded:
664 Length-adaptive average lagging for simultaneous speech translation. In Julia Ive and Ruiqing
665 Zhang (eds.), *Proceedings of the Third Workshop on Automatic Simultaneous Translation*, pp.
666 12–17, Online, July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.
667 autosimtrans-1.2. URL <https://aclanthology.org/2022.autosimtrans-1.2>.
- 668 Sara Papi, Matteo Negri, and Marco Turchi. Attention as a guide for simultaneous speech translation.
669 In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual
670 Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 13340–
671 13356, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/
672 v1/2023.acl-long.745. URL <https://aclanthology.org/2023.acl-long.745>.
- 673 Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic
674 evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Associa-
675 tion for Computational Linguistics*, pp. 311–318, Philadelphia, Pennsylvania, USA, July 2002.
676 Association for Computational Linguistics. doi: 10.3115/1073083.1073135. URL <https://aclanthology.org/P02-1040>.
- 677 Daniel S Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D Cubuk, and
678 Quoc V Le. Specaugment: A simple data augmentation method for automatic speech recognition.
679 *Interspeech 2019*, 2019.
- 680 Matt Post. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference
681 on Machine Translation: Research Papers*, pp. 186–191, Belgium, Brussels, October 2018. As-
682 sociation for Computational Linguistics. URL [https://www.aclweb.org/anthology/
683 W18-6319](https://www.aclweb.org/anthology/W18-6319).
- 684 Rohit Prabhavalkar, Kanishka Rao, Tara N. Sainath, Bo Li, Leif Johnson, and Navdeep Jaitly. A
685 Comparison of Sequence-to-Sequence Models for Speech Recognition. In *Proc. Interspeech
686 2017*, pp. 939–943, 2017. doi: 10.21437/Interspeech.2017-233.
- 687 Colin Raffel, Minh-Thang Luong, Peter J. Liu, Ron J. Weiss, and Douglas Eck. Online and
688 linear-time attention by enforcing monotonic alignments. In Doina Precup and Yee Whye Teh
689 (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of
690 *Proceedings of Machine Learning Research*, pp. 2837–2846. PMLR, 06–11 Aug 2017. URL
691 <https://proceedings.mlr.press/v70/raffell17a.html>.
- 692 Seamless Communication, Loïc Barrault, Yu-An Chung, Mariano Coria Meglioli, David Dale, Ning
693 Dong, Mark Duppenthaler, Paul-Ambroise Duquenne, Brian Ellis, Hady Elsahar, Justin Haaheim,
694 John Hoffman, Min-Jae Hwang, Hirofumi Inaguma, Christopher Klaiber, Iliia Kulikov, Pengwei
695 Li, Daniel Licht, Jean Maillard, Ruslan Mavlyutov, Alice Rakotoarison, Kaushik Ram Sadagopan,
696 Abinesh Ramakrishnan, Tuan Tran, Guillaume Wenzek, Yilin Yang, Ethan Ye, Ivan Evtimov,
697 Pierre Fernandez, Cynthia Gao, Prangthip Hansanti, Elahe Kalbassi, Amanda Kallet, Artyom
698 Kozhevnikov, Gabriel Mejia, Robin San Roman, Christophe Touret, Corinne Wong, Carleigh

- 702 Wood, Bokai Yu, Pierre Andrews, Can Balioglu, Peng-Jen Chen, Marta R. Costa-jussà, Maha
703 Elbayad, Hongyu Gong, Francisco Guzmán, Kevin Heffernan, Somya Jain, Justine Kao, Ann Lee,
704 Xutai Ma, Alex Mourachko, Benjamin Peloquin, Juan Pino, Sravya Popuri, Christophe Ropers,
705 Safiyyah Saleem, Holger Schwenk, Anna Sun, Paden Tomasello, Changhan Wang, Jeff Wang,
706 Skyler Wang, and Mary Williamson. Seamless: Multilingual expressive and streaming speech
707 translation. In *ArXiv*, 2023.
- 708 Frank Seide, Morrie Doulaty, Yangyang Shi, Yashesh Gaur, Junteng Jia, and Chunyang Wu. Speech
709 reallm – real-time streaming speech recognition with multimodal llms by teaching the flow of
710 time, 2024. URL <https://arxiv.org/abs/2406.09569>.
- 711
- 712 Chenze Shao and Yang Feng. Non-monotonic latent alignments for ctc-based non-
713 autoregressive machine translation. In S. Koyejo, S. Mohamed, A. Agarwal,
714 D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Process-*
715 *ing Systems*, volume 35, pp. 8159–8173. Curran Associates, Inc., 2022. URL
716 [https://proceedings.neurips.cc/paper_files/paper/2022/file/
717 35f805e65c77652efa731edc10c8e3a6-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/35f805e65c77652efa731edc10c8e3a6-Paper-Conference.pdf).
- 718 Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position represen-
719 tations. In Marilyn Walker, Heng Ji, and Amanda Stent (eds.), *Proceedings of the 2018 Con-*
720 *ference of the North American Chapter of the Association for Computational Linguistics: Hu-*
721 *man Language Technologies, Volume 2 (Short Papers)*, pp. 464–468, New Orleans, Louisiana,
722 June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-2074. URL
723 <https://aclanthology.org/N18-2074>.
- 724
- 725 Yangyang Shi, Yongqiang Wang, Chunyang Wu, Ching-Feng Yeh, Julian Chan, Frank Zhang, Duc
726 Le, and Mike Seltzer. Emformer: Efficient memory transformer based acoustic model for low
727 latency streaming speech recognition. In *ICASSP 2021 - 2021 IEEE International Conference*
728 *on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 6783–6787, 2021. doi: 10.1109/
729 ICASSP39728.2021.9414560.
- 730 Yun Tang, Anna Sun, Hirofumi Inaguma, Xinyue Chen, Ning Dong, Xutai Ma, Paden Tomasello,
731 and Juan Pino. Hybrid transducer and attention based encoder-decoder modeling for speech-
732 to-text tasks. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings*
733 *of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long*
734 *Papers)*, pp. 12441–12455, Toronto, Canada, July 2023. Association for Computational Linguis-
735 tics. doi: 10.18653/v1/2023.acl-long.695. URL [https://aclanthology.org/2023.
736 acl-long.695](https://aclanthology.org/2023.acl-long.695).
- 737 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
738 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von
739 Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Ad-*
740 *vances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.,
741 2017. URL [https://proceedings.neurips.cc/paper_files/paper/2017/
742 file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf).
- 743 Peidong Wang, Eric Sun, Jian Xue, Yu Wu, Long Zhou, Yashesh Gaur, Shujie Liu, and Jinyu Li.
744 LAMASSU: A Streaming Language-Agnostic Multilingual Speech Recognition and Translation
745 Model Using Neural Transducers. In *Proc. INTERSPEECH 2023*, pp. 57–61, 2023. doi: 10.
746 21437/Interspeech.2023-2004.
- 747
- 748 Zhichao Wang, Yuanzhe Chen, Xinsheng Wang, Lei Xie, and Yuping Wang. StreamVoice:
749 Streamable context-aware language modeling for real-time zero-shot voice conversion. In Lun-
750 Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting*
751 *of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 7328–7338,
752 Bangkok, Thailand, August 2024. Association for Computational Linguistics. URL [https://
753 aclanthology.org/2024.acl-long.396](https://aclanthology.org/2024.acl-long.396).
- 754 Chunyang Wu, Yongqiang Wang, Yangyang Shi, Ching-Feng Yeh, and Frank Zhang. Streaming
755 Transformer-Based Acoustic Models Using Self-Attention with Augmented Memory. In *Proc.*
Interspeech 2020, pp. 2132–2136, 2020. doi: 10.21437/Interspeech.2020-2079.

- 756 Jian Xue, Peidong Wang, Jinyu Li, Matt Post, and Yashesh Gaur. Large-Scale Streaming End-to-End
757 Speech Translation with Neural Transducers. In *Proc. Interspeech 2022*, pp. 3263–3267, 2022.
758 doi: 10.21437/Interspeech.2022-10953.
- 759
760 Kingshan Zeng, Liangyou Li, and Qun Liu. RealTrans: End-to-end simultaneous speech transla-
761 tion with convolutional weighted-shrinking transformer. In Chengqing Zong, Fei Xia, Wenjie
762 Li, and Roberto Navigli (eds.), *Findings of the Association for Computational Linguistics: ACL-*
763 *IJCNLP 2021*, pp. 2461–2474, Online, August 2021. Association for Computational Linguis-
764 tics. doi: 10.18653/v1/2021.findings-acl.218. URL [https://aclanthology.org/2021.
765 findings-acl.218](https://aclanthology.org/2021.findings-acl.218).
- 766 Qian Zhang, Han Lu, Hasim Sak, Anshuman Tripathi, Erik McDermott, Stephen Koo, and Shankar
767 Kumar. Transformer transducer: A streamable speech recognition model with transformer en-
768 coders and rnn-t loss. In *ICASSP 2020 - 2020 IEEE International Conference on Acoustics,
769 Speech and Signal Processing (ICASSP)*, pp. 7829–7833, 2020. doi: 10.1109/ICASSP40776.
770 2020.9053896.
- 771 Ruiqing Zhang, Zhongjun He, Hua Wu, and Haifeng Wang. Learning adaptive segmentation pol-
772 icy for end-to-end simultaneous translation. In Smaranda Muresan, Preslav Nakov, and Aline
773 Villavicencio (eds.), *Proceedings of the 60th Annual Meeting of the Association for Computa-*
774 *tional Linguistics (Volume 1: Long Papers)*, pp. 7862–7874, Dublin, Ireland, May 2022. As-
775 sociation for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.542. URL [https:
776 //aclanthology.org/2022.acl-long.542](https://aclanthology.org/2022.acl-long.542).
- 777 Shaolei Zhang and Yang Feng. Unified segment-to-segment framework for simul-
778 taneous sequence generation. In A. Oh, T. Naumann, A. Globerson, K. Saenko,
779 M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing*
780 *Systems*, volume 36, pp. 45235–45258. Curran Associates, Inc., 2023. URL
781 [https://proceedings.neurips.cc/paper_files/paper/2023/file/
782 8df705957a5262de3cb37ba9f1fb96f3-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/8df705957a5262de3cb37ba9f1fb96f3-Paper-Conference.pdf).
- 783 Shaolei Zhang, Qingkai Fang, Shoutao Guo, Zhengrui Ma, Min Zhang, and Yang Feng. Stream-
784 speech: Simultaneous speech-to-speech translation with multi-task learning. In *Proceedings of*
785 *the 62th Annual Meeting of the Association for Computational Linguistics (Long Papers)*. Asso-
786 ciation for Computational Linguistics, 2024.
- 787
788 Jinzheng Zhao, Niko Moritz, Egor Lakomkin, Ruiming Xie, Zhiping Xiu, Katerina Zmolikova,
789 Zeeshan Ahmed, Yashesh Gaur, Duc Le, and Christian Fuegen. Textless streaming speech-to-
790 speech translation using semantic speech tokens. *arXiv preprint arXiv:2410.03298*, 2024.
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

A IMPLEMENTATION DETAILS

Pre-processing The input speech is represented as 80-dimensional log mel-filterbank coefficients computed every 10ms with a 25ms window. Global channel mean and variance normalization is applied to the input speech. During training, SpecAugment (Park et al., 2019) data augmentation with the LB policy is additionally employed. We use SentencePiece (Kudo & Richardson, 2018) to generate a unigram vocabulary of size 10000 for the source and target text jointly. Sequence-level knowledge distillation (Kim & Rush, 2016) is applied for fair comparison (Liu et al., 2021). For speech-to-speech experiments, we resample the source audio to 16kHz and apply identical pre-processing steps as those used in speech-to-text experiments. For the target speech, we also downsample the audio and extract discrete units utilizing the publicly available pre-trained mHuBERT model and K-means quantizer.⁴ No training data manipulation is applied in speech-to-speech experiments.

Training Details Considering that training MonoAttn-Transducer involves two critical processes: inferring the posterior alignment and estimating the context vector, instability in either step can lead to training failure. Therefore, we introduce a curriculum learning strategy for MonoAttn-Transducer. We first pretrain the model in an offline setting. In pretraining, all predictor states can attend to the complete source input, and the model is trained as an offline Transducer. This pretraining phase allows the monotonic attention module to warm up by learning full-sentence attention, thereby enhancing its stability during subsequent adaptation to a streaming scenario. In finetuning, we apply Algorithm 1 to adjust MonoAttn-Transducer with various chunk size configurations. During both training phases, we set the dropout rate to 0.1, weight decay to 0.01, and clip gradient norms exceeding 5.0. The dropout rates for activation and attention are both set to 0.1. The pretraining spans 50k updates with a batch size of 160k tokens. The learning rate gradually warms up to 5e-4 within 4k steps. Finetuning involves training for 20k updates and other hyper-parameters remain consistent. Throughout the training, we optimize models using the Adam optimizer (Kingma & Ba, 2015). Automatic mixed precision training is applied. It takes approximately one day to pretrain in an offline setting and another day for streaming adaptation on a server with 4 Nvidia L40 GPUs.

B BASELINES

We compare our proposed MonoAttn-Transducer with the following state-of-the-art open-source approaches (without using pretrained encoder or any data augmentation method for fair comparison).

AED-BASED MODELS

Wait- k (Ma et al., 2020c): It executes wait- k policy (Ma et al., 2019) by setting the pre-decision window size to 280 ms.

RealTrans (Zeng et al., 2021): It detects word number in the streaming speech by counting blanks in CTC transcription and applies wait- k -stride- n strategy accordingly.

MU-ST (Zhang et al., 2022): It trains an external segmentation model, which is then utilized to detect meaningful units for guiding generation.

Seg2Seg (Zhang & Feng, 2023): It alternates between waiting for a source segment and generating a target segment in an autoregressive manner.

EDAtt (Papi et al., 2023): It calculates the attention scores towards the latest received frames of speech, serving as guidance for an offline-trained translation model during simultaneous inference.

CTC-BASED MODELS

NAST (Ma et al., 2024): It introduces a streaming generation model with fast computation speed by leveraging a non-autoregressive transformer and CTC decoding (Graves et al., 2006).

⁴https://github.com/facebookresearch/fairseq/blob/main/examples/speech_to_speech/docs/textless_s2st_real_data.md

864 TRANSDUCER-BASED MODELS

865
866 **Transducer:** It adopts the standard Transducer framework (Graves, 2012) and utilizes Transformer
867 as its backend network (Zhang et al., 2020).868 **CAAT** (Liu et al., 2021): It incorporates a cross-attention module within Transducer’s joiner to
869 alleviate its strong monotonic constraint.

871 C NUMERICAL RESULTS

872
873 In addition to Average Lagging (AL; Ma et al., 2020c), we also incorporate Average Proportion (AP;
874 Cho & Esipova, 2016), Differentiable Average Lagging (DAL; Arivazhagan et al., 2019) and Length
875 Adaptive Average Lagging (LAAL; Papi et al., 2022) as metrics to evaluate the latency. AL, DAL and
876 LAAL are all reported with milliseconds. The trade-off between latency and translation quality is
877 attained by adjusting the chunk size C . The offline results are obtained by setting the chunk size to
878 be longer than any utterance in the dataset ($C = \infty$). We use SimulEval v1.1.4 for evaluation in
879 all the experiments. The numerical results of MonoAttn-Transducer are presented in Table 5 and 6.
880 A comparison of the *computation-aware* latency metrics for AL and LAAL between the Transducer
881 and MonoAttn-Transducer models is presented in Table 7.882
883 Table 5: Numerical results of MonoAttn-Transducer on MuST-C English to German dataset.884
885

<i>MonoAttn-Transducer on En→De</i>					
$C(ms)$	AP	AL	DAL	LAAL	BLEU
320	0.67	1215	1497	1317	20.22
640	0.77	1470	1872	1582	22.47
960	0.83	1860	2309	1957	22.94
1280	0.86	2215	2719	2305	23.74
∞	-	-	-	-	24.42

892
893 Table 6: Numerical results of MonoAttn-Transducer on MuST-C English to Spanish dataset.894
895

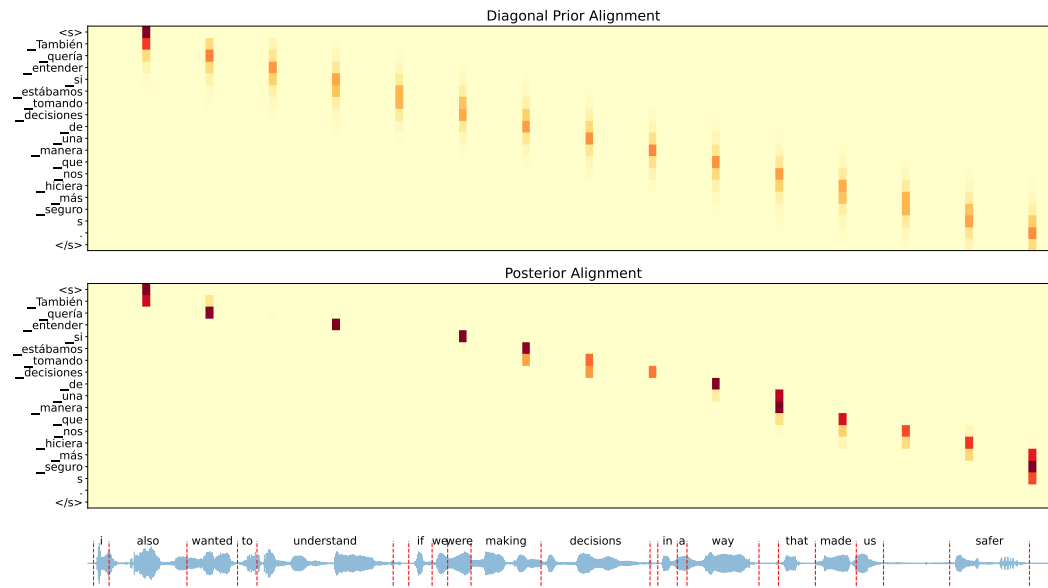
<i>MonoAttn-Transducer on En→Es</i>					
$C(ms)$	AP	AL	DAL	LAAL	BLEU
320	0.74	997	1534	1230	24.72
640	0.81	1239	1854	1475	26.74
960	0.88	1606	2304	1837	27.05
1280	0.93	1991	2725	2204	27.41
∞	-	-	-	-	27.48

903
904 Table 7: Comparison of MonoAttn-Transducer and Transducer across various chunk size settings
905 on MuST-C English to German and English to Spanish datasets.906
907

		<i>En-Es</i>				<i>En-De</i>				
		Chunk Size (ms)	320	640	960	1280	320	640	960	1280
Transducer	AL (ms, \downarrow)		886	1193	1591	1997	1126	1434	1830	2215
	AL_CA (ms, \downarrow)		1121	1330	1699	2085	1323	1551	1920	2296
	LAAL (ms, \downarrow)		1168	1466	1847	2220	1258	1563	1942	2312
	LAAL_CA (ms, \downarrow)		1381	1589	1944	2300	1444	1673	2028	2389
MonoAttn-Transducer	AL (ms, \downarrow)		997	1239	1606	1991	1215	1470	1860	2215
	AL_CA (ms, \downarrow)		1239	1385	1724	2089	1407	1596	1964	2301
	LAAL (ms, \downarrow)		1230	1475	1837	2204	1317	1582	1957	2305
	LAAL_CA (ms, \downarrow)		1453	1607	1945	2295	1501	1702	2056	2387

Table 8: Comparison of results reported in Tang et al. (2023) and our work on MuST-C English to German dataset.

	Chunk Size (ms)	160	320	480	640
Transducer (Tang et al., 2023)	BLEU (\uparrow)	20.76	21.80	22.52	23.32
	AL (ms, \downarrow)	1282	1252	1306	1498
TAED (Tang et al., 2023)	BLEU (\uparrow)	21.57	22.63	23.48	23.47
	AL (ms, \downarrow)	1263	1354	1369	1903
	Chunk Size (ms)	320	640	960	1280
Transducer (Our implementation)	BLEU (\uparrow)	19.99	22.10	22.20	22.96
	AL (ms, \downarrow)	1126	1434	1830	2215
MonoAttn-Transducer	BLEU (\uparrow)	20.22	22.47	22.94	23.74
	AL (ms, \downarrow)	1215	1470	1860	2215

Figure 4: Chunk size in this example is set to 320ms. (*Diagonal Prior*)

D VISUALIZATION

In this section, we present more examples of *diagonal prior* and its posterior from training corpus. Additionally, we also provide examples of *uniform prior* and its posterior for comparison. We have observed that, even with significant differences in the prior distribution, the posterior remains fairly robust when the chunk size is constant. The vertical axis represents the target subword sequence and the horizontal axis represents the speech waveform. Darker areas indicate higher alignment probabilities. We use Montreal Forced Alignment tools (McAuliffe et al., 2017) to obtain speech-transcription alignments for illustration.

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1000

1001

1002

1003

1004

1005

1006

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

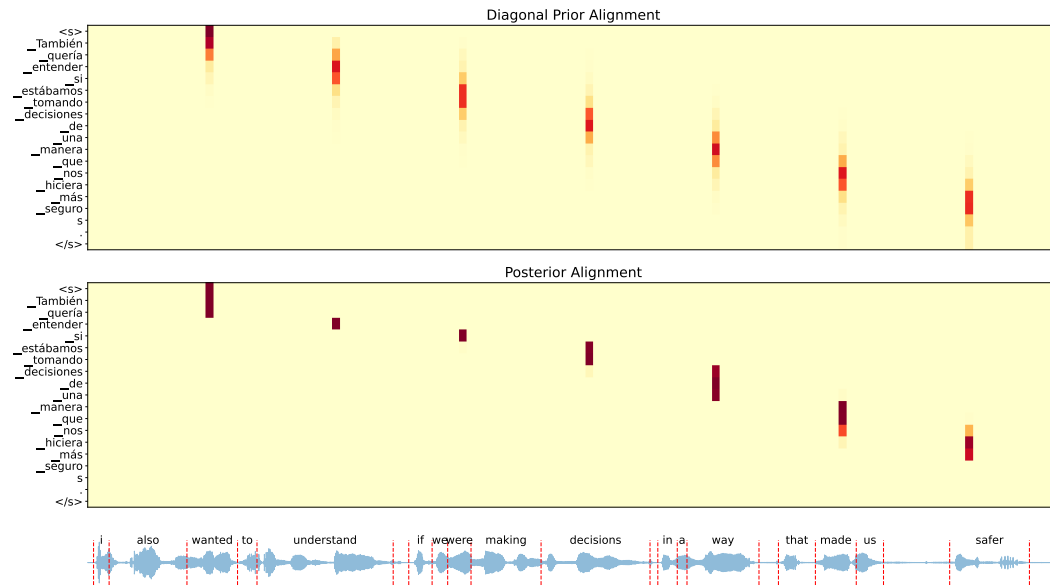


Figure 5: Chunk size in this example is set to 640ms. (*Diagonal Prior*)

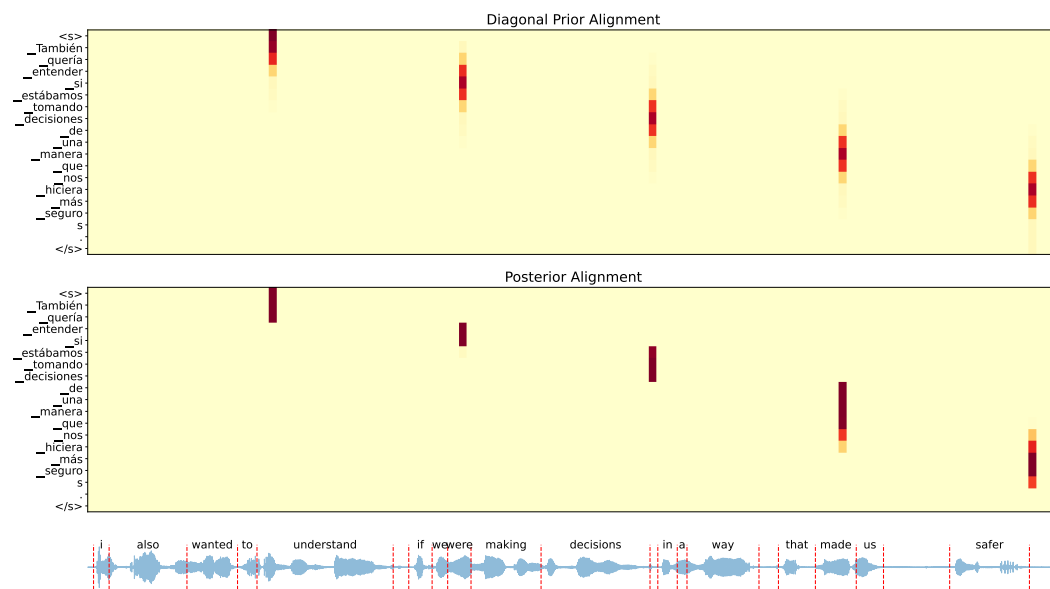


Figure 6: Chunk size in this example is set to 960ms. (*Diagonal Prior*)

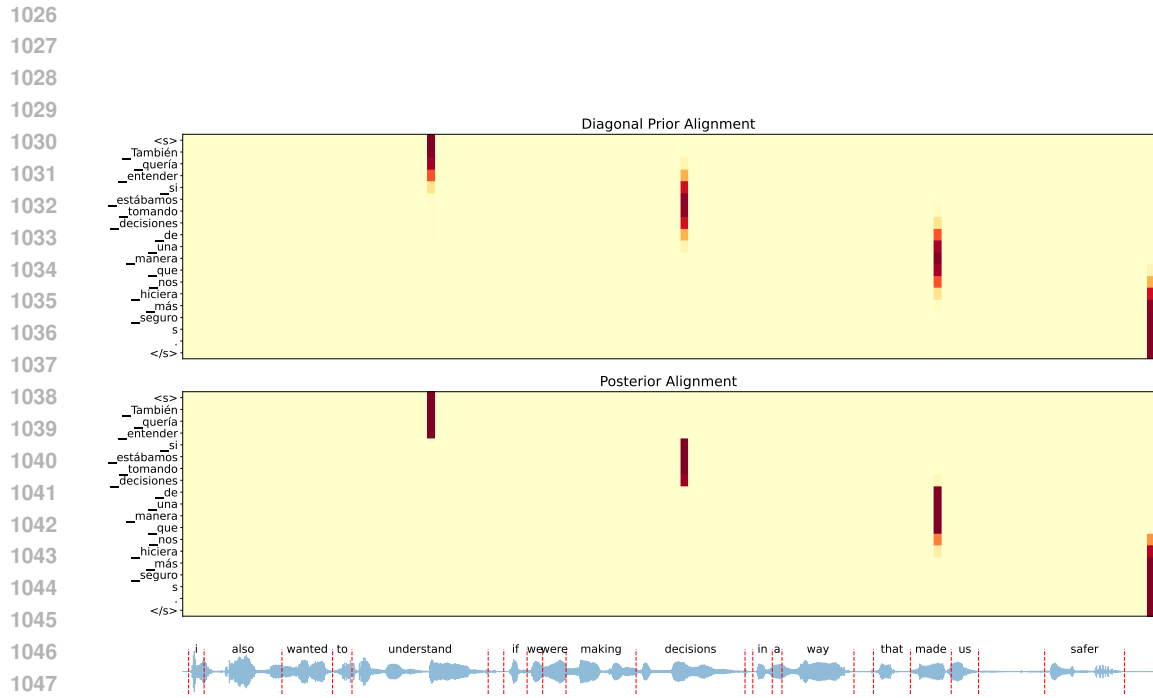


Figure 7: Chunk size in this example is set to 1280ms. (*Diagonal Prior*)

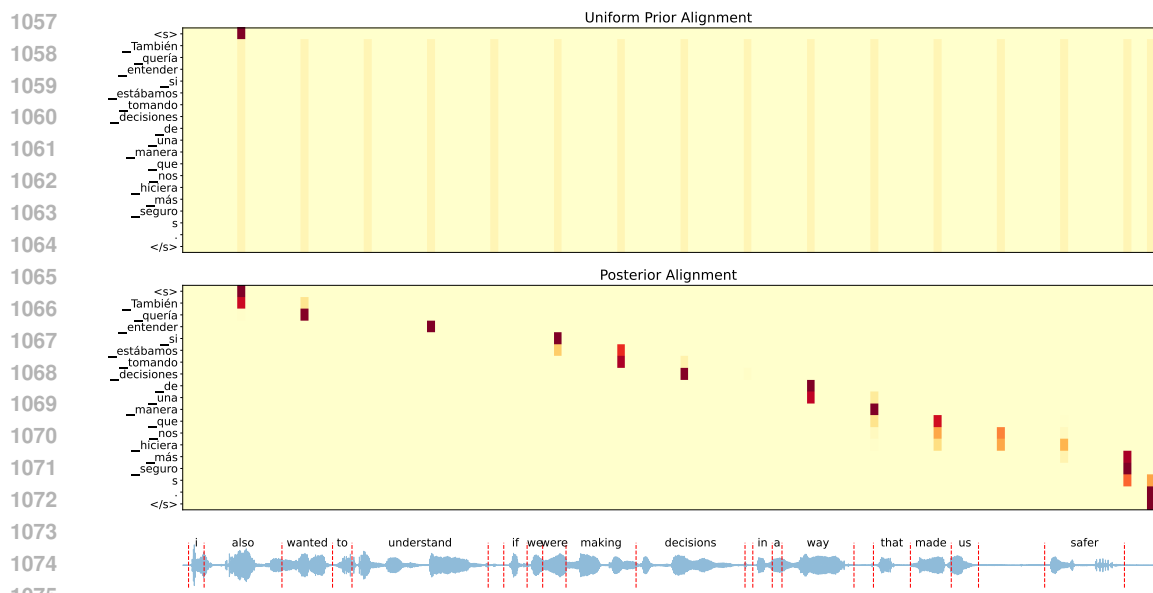
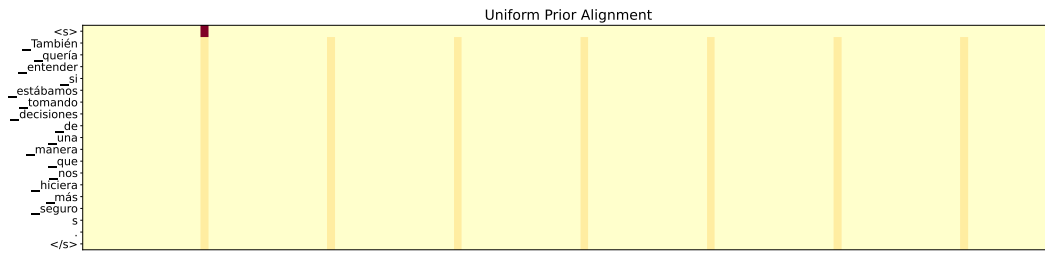
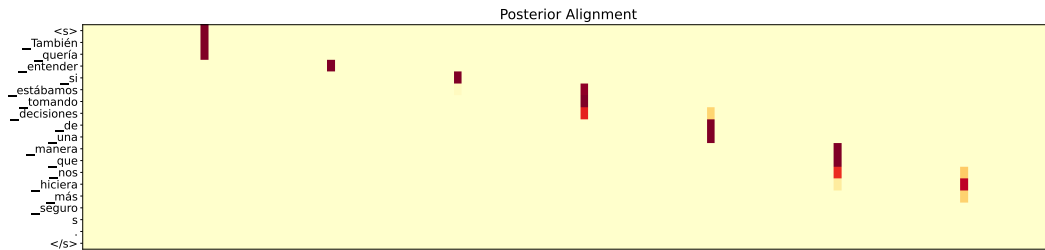


Figure 8: Chunk size in this example is set to 320ms. (*Uniform Prior*)

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091



1092
1093
1094
1095
1096
1097
1098
1099



1100
1101

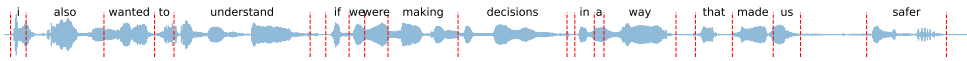
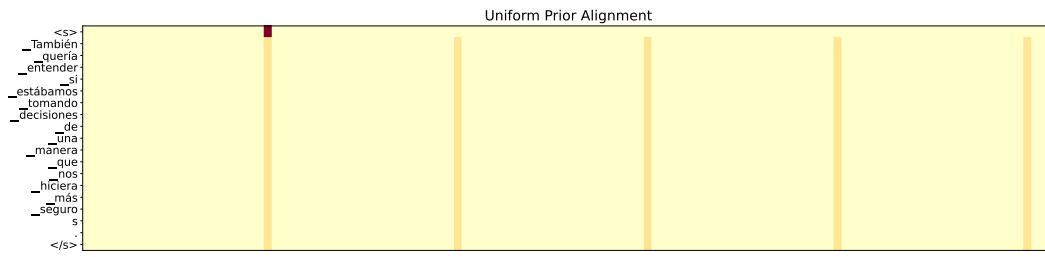


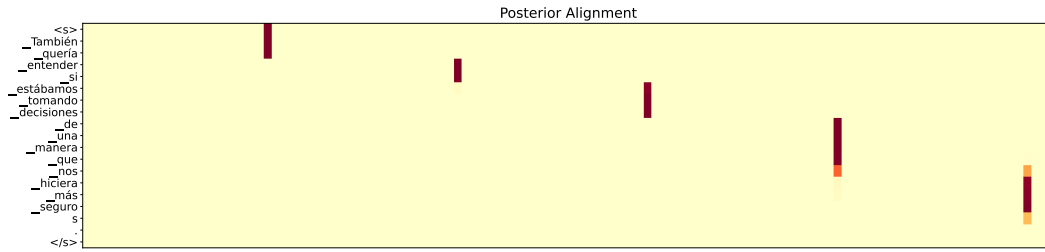
Figure 9: Chunk size in this example is set to 640ms. (*Uniform Prior*)

1102
1103
1104
1105
1106
1107
1108
1109
1110

1111
1112
1113
1114
1115
1116
1117
1118



1119
1120
1121
1122
1123
1124
1125
1126



1127
1128

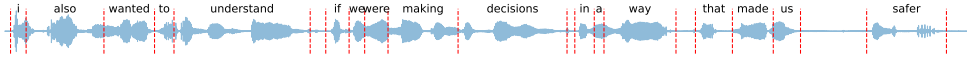


Figure 10: Chunk size in this example is set to 960ms. (*Uniform Prior*)

1129
1130
1131
1132
1133

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

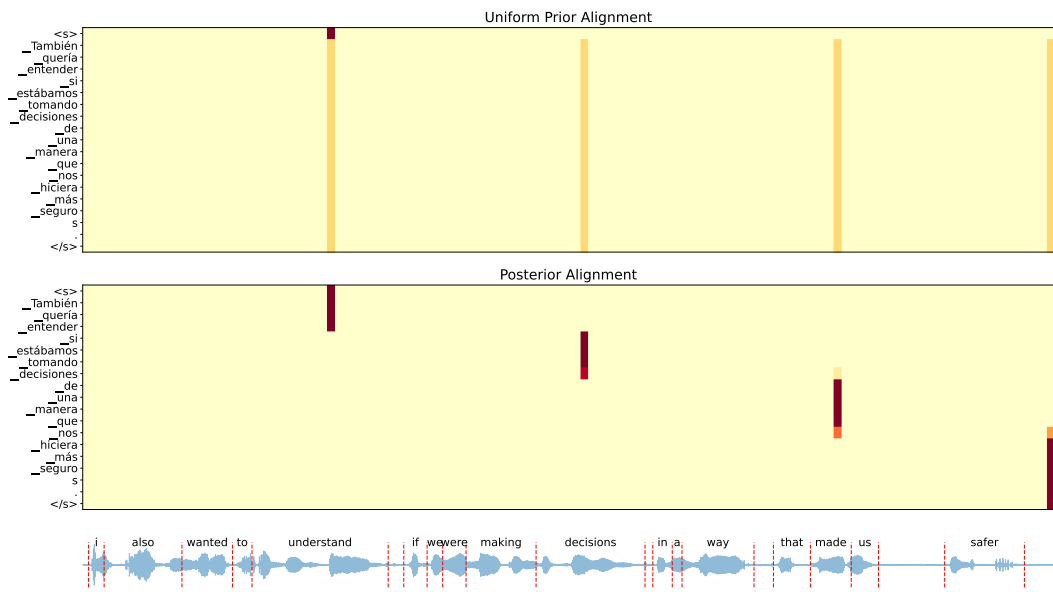


Figure 11: Chunk size in this example is set to 1280ms. (*Uniform Prior*)