

IMPROVING FAST-SLOW ENCODER BASED TRANSDUCER WITH STREAMING DELIBERATION

Ke Li, Jay Mahadeokar, Jinxi Guo, Yangyang Shi, Gil Keren, Ozlem Kalinli, Michael L. Seltzer, Duc Le

Meta AI, USA

ABSTRACT

This paper introduces a fast-slow encoder based transducer with streaming deliberation for end-to-end automatic speech recognition. We aim to improve the recognition accuracy of the fast-slow encoder based transducer while keeping its latency low by integrating a streaming deliberation model. Specifically, the deliberation model leverages partial hypotheses from the streaming fast encoder and implicitly learns to correct recognition errors. We modify the parallel beam search algorithm for fast-slow encoder based transducer to be efficient and compatible with the deliberation model. In addition, the deliberation model is designed to process streaming data. To further improve the deliberation performance, a simple text augmentation approach is explored. We also compare LSTM and Conformer models for encoding partial hypotheses. Experiments on Librispeech and in-house data show relative WER reductions (WERRs) from 3% to 5% with a slight increase in model size and negligible extra token emission latency compared with fast-slow encoder based transducer. Compared with vanilla neural transducers, the proposed deliberation model together with fast-slow encoder based transducer obtains relative 10-11% WERRs on Librispeech and around relative 6% WERR on in-house data with smaller emission delays.

Index Terms— Fast-slow encoder-based transducer, Deliberation model, Parallel beam search, RNN-T, Conformer

1. INTRODUCTION

Low latency is a general requirement for voice assistants for good user experience. To this end, the streaming capability of an automatic speech recognition (ASR) system is a critical consideration. The recurrent neural network transducer (RNN-T) [1–4], which intrinsically supports streaming transcription, is a commonly adopted model for end-to-end (E2E) ASR [1–9].

Streaming RNN-Ts limit the lookahead context access of the input audio for latency control. Various two-pass approaches have been introduced to address the accuracy loss due to limited context [10–12]. For example, an attention-based encoder-decoder model was used to rescore N-best hypotheses generated by an RNN-T in the 1st-pass [10]. The deliberation model [11] subsequently extended the previous work by adding a spelling correction component [13]. Another way to improve the performance of streaming RNN-Ts is to introduce a non-causal encoder with larger context. Narayanan et al. [14] proposed a cascaded encoder-based RNN-T with a causal/streaming encoder and a non-causal/non-streaming encoder. Li et al. [15] applied two-pass beam search with cascaded encoders, where the 1st-pass uses a streaming encoder and the 2nd-pass uses a non-streaming encoder with left and right context.

The non-streaming encoder in the cascaded architecture, however, introduces non-negligible latency. For latency-constrained use case, Mahadeokar et al. [16] proposed a streaming fast-slow encoder

based transducer design with both encoders having limited future lookahead. This streaming design achieves substantial accuracy improvement by the slow encoder, which takes multiple segments' output of the fast encoder as the input. Mahadeokar et al. [16] also proposed a novel streaming parallel beam search for the fast-slow encoders where the search space for fast and slow encoders is shared.

However, the accuracy of the fast-slow encoder based transducer [16] is still compromised since the lookahead context of the streaming slow encoder is limited. In order to improve its accuracy while keeping the low latency, we integrate a streaming deliberation model into fast-slow encoder based transducer that leverages partial hypotheses decoded from the fast encoder. Specifically, the partial hypotheses are passed through a text encoder. The resulting text embeddings are combined with the acoustic embeddings from the slow encoder by a merge module. The combined embeddings are then fed to the joiner. To allow the deliberation model to process streaming data, we limit the maximum length of partial hypotheses for the text encoder. During training, random masks are applied to partial hypotheses as a simple augmentation approach. Unlike [17] where an extra joiner is used, we use a shared joiner initialized by the joiner in a fast-slow encoder based transducer. This simplifies the model design and makes training converge faster. We integrate the deliberation model into the parallel beam search algorithm. During inference, partial hypotheses are only generated by beam search from the fast encoder that is immediately before the call of the slow encoder beam search. Furthermore, we compare different text encoder architectures in terms of accuracy and emission delay, and discuss limitations of the deliberation approach.

2. METHODS

2.1. Fast-Slow Encoder Based Transducer

We first introduce the baseline model, the streaming fast-slow encoder based transducer [16]. Compared to RNN-Ts with a single encoder, fast-slow encoder based transducer improves recognition accuracy by leveraging more acoustic information through the slow encoder, which effectively has K times acoustic context of the fast encoder. Different from the causal and non-causal cascaded encoders approach [14], both fast and slow encoders in [16] are streaming and has limited lookahead context for latency control. A novel parallel beam search algorithm was developed for fast-slow encoder based transducer, where the slow encoder can update hypotheses decoded from the fast encoder and the search space for the two encoders is shared. In summary, the fast-slow encoder based transducer achieves much better recognition accuracy without much increase in token emission delay compared to a baseline streaming RNN-T with similar architecture. In this work, we build the proposed streaming deliberation model on top of the fast-slow encoder based transducer, where the fast encoder is used to generate partial hypotheses for the streaming deliberation model.

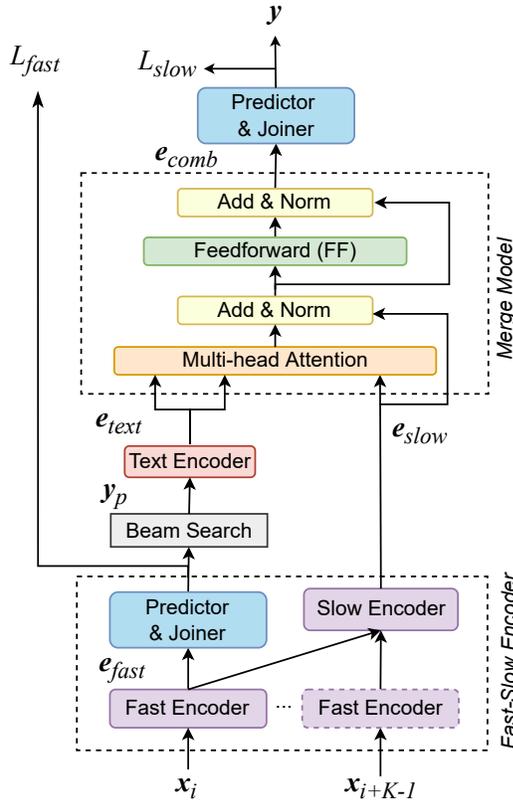


Fig. 1. Fast-slow encoder based transducer with the streaming deliberation model.

2.2. Streaming Deliberation Model for Fast-Slow Encoder Based Transducer

In this section, we introduce the fast-slow encoder based transducer with streaming deliberation including model details, training approach, and a simple text augmentation approach.

2.2.1. Model

The deliberation model improves the recognition accuracy of fast-slow encoder based transducer by leveraging partial decoded hypotheses from the fast encoder. The model is illustrated in Fig. 1. The fast-slow encoder based transducer is shown at the bottom of Fig. 1. To support streaming processing, input acoustic frames are segmented into chunks and each chunk contains a certain number of audio frames. Each processed audio chunk \mathbf{x}_i is concatenated with a small lookahead context and then fed into the fast encoder. The slow encoder is cascaded with the fast encoder and takes around K outputs of the fast encoder as the input, e.g., K can be two or more, as shown in Fig. 1. Both the fast and slow encoders take a lookahead context size 1, which is effectively 40ms.

We perform decoding on training data with the fast encoder to generate 1-best or N -best partial hypotheses. Let us denote 1-best hypothesis as \mathbf{y}_p . It is then encoded by a text encoder, which consists of a lookup embedding matrix and a neural network-based text encoder. The text embeddings \mathbf{e}_{text} from the text encoder are combined with the acoustic embedding \mathbf{e}_{slow} from the slow encoder by the ‘Merge Model’ in Fig. 1. The ‘Merge Model’ consists of multi-

ple blocks, where each block has a multi-head attention module and a feedforward (FF) network. The acoustic embedding \mathbf{e}_{slow} is the query, and the text embeddings \mathbf{e}_{text} are key and value in the attention module. The residual connection adds a summarized text embedding with each acoustic embedding. The combined embedding from the ‘Merge Model’ denoted as \mathbf{e}_{comb} is then fed to the joiner along with the output of the predictor, which is omitted in Fig. 1 for simplicity. The predictor and joiner of the deliberation model share their weights with the predictor and joiner that connects to the fast encoder, respectively. To further reduce the model size, we also share the token embedding matrices between the text encoder and the predictor.

We compare the use of LSTM and Conformer architectures for the text encoder. To support streaming processing, the LSTM network has to be unidirectional. For Conformer [18], we set a limit on the maximum number of tokens a partial hypothesis can contain; we found that this truncation with maximum length 20 has negligible impact on the word error rate (WER). Our implementation of the streaming deliberation model is different from [17] considering the slow encoder beam search updates the output of the fast encoder so partial decoded hypotheses from the fast encoder vary over time. The partial hypotheses are passed through the text encoder every time before the beam search call over the slow encoder. Since the number of beam search calls over the slow encoder are much fewer than the number of beam search calls over the fast encoder, the extra computation cost of the deliberation model is not much.

2.2.2. Joint Training

The training objective of the proposed fast-slow encoders with a streaming deliberation module is

$$L = L_{slow} + \lambda L_{fast} \quad (1)$$

where $0 < \lambda < 1$. Both L_{fast} and L_{slow} use alignment restricted RNN-T loss [4]. As shown in Fig. 1, L_{fast} is computed from the fast encoder and L_{slow} is computed from the slow encoder with deliberation module. Different from previous work [14] where the causal encoder and the non-causal encoder randomly use different training samples within a minibatch, both fast and slow encoders in our model use all the training data.

The training procedure has two steps:

- Train the fast-slow encoder based transducer as in [4].
- Jointly train the fast-slow encoder based transducer initialized from the first step with the deliberation model.

In the second training step, the fast-slow encoder based transducer is jointly optimized with the deliberation model. We experimented with freezing all or part of fast-slow encoder based transducer model parameters but neither outperformed joint training.

Inspired by [19], we experiment with a simple text augmentation method to increase the diversity of the partial hypotheses for deliberation. We randomly mask out word piece tokens in partial hypotheses with a small probability during training. The mask token we use is the blank token, so there is no need to change the vocabulary. This is different from alignment augmentation [20] where the blank token is already included in decoded hypotheses thus a new symbol that represents the mask token is required. Similar to [19], this random masking trick is applied only in training as an augmentation approach, thus no change is required in inference.

2.3. Parallel Beam Search with Deliberation Model

We modify the original parallel beam search algorithm for the fast-slow encoder based transducer [16] to support the streaming deliberation model. The modified algorithm is described below (modifications are in purple).

Algorithm 1 Parallel beam search for fast-slow encoders with streaming deliberation.

```

 $T^f = \text{fastSegmentSize}()$ 
 $T^s = \text{slowSegmentSize}()$ 
 $N^f = \text{fastBeamSize}()$ 
 $N^s = \text{slowBeamSize}()$ 
 $B^f \leftarrow \emptyset; B^s \leftarrow \emptyset$ 
 $(H^f, H^s) \leftarrow \text{initModelState}()$ 
 $\Gamma \leftarrow \text{initSearchSpace}()$ 
 $I^s \leftarrow \emptyset$ 
for  $t = T^f$  to  $T$  by  $T^f$  do
   $(I^f, R^f) = \text{getFeatures}(t, t - T^f)$ 
   $(e_f, H^f) = \text{encodeFast}(I^f, (I^f, R^f, H^f))$ 
   $B^f \leftarrow \text{beamSearch}(e_f, B^f, N^f, \Gamma, H^f)$ 
  if  $t \bmod T^s = T^s - 1$  then
     $y_p = \text{getBestHypo}(B^f)$ 
  end if
   $I^s \leftarrow \text{concat}(I^s, e_f)$ 
  if  $t \bmod T^s = 0$  or  $t = T$  then
     $(e^s, H^s) = \text{encodeSlow}(I^s, (I^s, R^s, H^s))$ 
     $e^{\text{comb}} = \text{encodeDeliberation}(e_s, y_p)$ 
     $B^s \leftarrow \text{beamSearch}(e^{\text{comb}}, B^s, N^s, \Gamma, H^s)$ 
     $I^s \leftarrow \emptyset$ 
     $B^f \leftarrow B^s$ 
  end if
end for
return  $y$  with highest  $\log Pr(y)/|y|$  in  $B^s$ 

```

Note that ‘f’ denotes ‘fast’ and ‘s’ denotes ‘slow’ in the above algorithm. Let us assume Γ represents the shared decoding space for fast and slow encoders. T^f and T^s denote chunk or segment size of fast and slow encoders and T denotes acoustic sequence length. N^f and N^s denote beam sizes for decoding from the fast encoder and the slow encoder respectively. B^f and B^s represent the N-best hypotheses generated using fast-slow encoders, respectively. H^f and H^s denote decoding states for fast and slow encoders. The 1-best partial hypothesis from the fast encoder is denoted as y_p .

The key process of the algorithm contains interval calls of beam search over fast and slow encoders. After initialization, we iterate over acoustic frames from time step 0 to T in the interval of T^f and run beam search with the fast encoder. The output acoustic embedding e^f from the fast encoder is accumulated until we need to perform beam search with the slow encoder. The accumulated embedding represented by I^s is used as input for the slow encoder. After the beam search call of the slow encoder, we update B^f with B^s and run the next fast encoder beam search on the next chunk of acoustic frames until the end.

The highlighted part in purple shows how deliberation is integrated into the algorithm. When we process the chunk of acoustic frames just before the next beam search call over the slow encoder, we trace the 1-best partial hypothesis y_p from the current beam search over the fast encoder. The `encodeDeliberation` operation then encodes y_p and combines it with acoustic embedding e^s from the slow encoder by the merge model. The resulting embedding e^{comb} is then used for beam search to generate B^f . Note

that `encodeDeliberation` call happens infrequently, only after `encodeSlow` is called. This saves computation from the deliberation model.

3. EXPERIMENTAL SETUP

3.1. Datasets

3.1.1. Librispeech

The Librispeech [21] corpus contains 960 hours of labeled speech. 80-dimensional filter bank features are extracted from a 25 ms window with a stride of 10 ms. We apply SpecAugment [22] with mask parameter $F = 27$, ten time masks with maximum time-mask ratio $p_s = 0.05$, and speed perturbation.

3.1.2. Large-Scale In-House Data

Our in-house training set combines two sources. The first consists of 20K hours of English video data publicly shared by Facebook users; all videos are completely de-identified before transcription. The second contains 20K hours of manually transcribed de-identified English data with no user-identifiable information in the voice assistant domain. All utterances are morphed when researchers manually access them to further de-identify the speaker. Note that the data is not morphed during training. We further augment the data with speed perturbation, simulated room impulse response, and background noise, resulting in 83M utterances (145K hours).

We evaluate our models on two in-house test sets:

VA1 – 10.2K hand-transcribed de-identified short-form utterances (less than five words on average) in the voice assistant domain, collected from internal volunteers. The participants consist of households that have agreed to have their voice activity reviewed and analyzed.

VA2 – 44.2K hand-transcribed de-identified short-form utterances in the voice assistant domain, collected by a third-party data vendor.

3.2. Models

Our baseline RNN-T models use 20 layers of Emformer [23] as the encoder. The joiner is a 1-layer feed-forward network, and the predictor is a 3-layer LSTM. The word piece vocabulary size is 5001 for Librispeech and 4096 for in-house data. The total number of parameters of the baseline RNN-T model is approximately 79M on Librispeech and 78M on in-house data. Baseline fast-slow encoder based transducer consists of 15 fast encoder layers and 5 slow encoder layers. Each encoder layer uses the Emformer architecture. We choose $\lambda = 0.5$ for the fast encoder loss (Equation 1). We experiment with both LSTM and Conformer for text encoders. For Conformers, we limit the total token size of each partial hypothesis up to 20. The Merge Model consists of 1-layer multi-head attention with one head and a feed-forward network.

We evaluate recognition accuracy using WER and measure token emission latency by tracking average, P95 and P99 emission delays. Emission delay is defined as the time from when the token is spoken to when the transcript of the token is emitted [4].

Parameter setups to train baseline RNN-T and fast-slow encoder based transducer follow [16]. For training the streaming deliberation models, we use a 100 times smaller learning rate and half the number of total epochs compared to the baseline models. To improve training efficiency, we set the beam size to 1 when decoding training data. Compared to larger beam sizes, e.g. 5 or 10, beam size 1

significantly accelerates training without degrading model accuracy. During inference, we use beam size 10 for all models.

4. EXPERIMENTAL RESULTS

4.1. WERs and Latency on Librispeech

We compare the WER and emission delay (ED) of the proposed model and baseline RNN-T and fast-slow encoder based transducer on Librispeech. The total extra parameters from the deliberation module with a 3-layer Conformer text encoder is around 16M. We use context size 160ms for the fast encoder and 800ms for the slow encoder. As shown in Table 1, the deliberation model with random mask probability 0.1 achieves 5% and 3% WERR on test-clean and test-other, respectively, compared to fast-slow encoder based transducer. It slightly increases average emission delay while having no effect on P95 and P99 compared to fast-slow encoders. The relative WER reductions are 10-11% compared to RNN-Ts. As for random masking, small yet consistent gains are achieved by using small probabilities such as 0.05 and 0.1. Using a masking probability greater than 0.2 does not show improvement.

Table 1. WERs (%) and ED (ms) from baseline RNN-Ts and proposed fast-slow encoder based transducer on Librispeech.

model	context	WER		ED		
		test-clean	test-other	avg	P95	P99
RNN-T	160	3.54	8.85	344	480	560
fast-slow		3.32	8.17	335	480	600
+ deliberation	160/800	3.20	7.98	337	480	600
+ mask p = 0.1		3.15	7.90	337	480	600

4.2. Effect of Text Encoder Architecture

Text encoder is a key component in the deliberation model and can affect both accuracy and latency. We experiment with both LSTM and Conformer architectures in 1-layer and 3-layer setups. All deliberation models are trained with a random masking probability of 0.1. Table 2 presents the comparison results on Librispeech. For 1-layer text encoders, LSTM performs slightly better than Conformer. This is expected since usually deeper Conformer architectures are more powerful. For the 3-layer setup, Conformer performs slightly better. There is no significant difference in ED for both architectures.

Table 2. WERs (%) and ED (ms) from LSTM and Conformer text encoders of the deliberation model on Librispeech.

model	context	WER		ED	
		test-clean	test-other	avg	P99
RNN-T	160	3.54	8.85	344	560
fast-slow	160/800	3.32	8.17	335	600
+ delib (LSTM 1-L)	160/800	3.17	7.96	337	600
+ delib (Conformer 1-L)		3.25	8.01	337	600
+ delib (LSTM 3-L)	160/800	3.20	7.93	337	600
+ delib (Conformer 3-L)		3.15	7.90	337	600

4.3. Limitation of Deliberation

The key capacity of deliberation model is its potential correction ability learned by observing hypotheses with errors in training. Therefore, one key limitation of this approach is that it may not help much for short utterances without enough errors. To test this hypothesis, we conduct an analysis on the Librispeech dataset. We split each test set into two parts, one containing utterances shorter than 3s and the other with utterances longer than 3s. WERs in Table 3 show improvements are mainly from longer utterances as hypothesized. This indicates the deliberation model is more powerful for dictation application.

Table 3. WERs (%) from test sets with different duration on Librispeech.

model	test-clean		test-other	
	>3s	≤3s	>3s	≤3s
fast-slow	3.3	4.3	8.0	10.3
+ deliberation	3.1	4.4	7.7	10.3

4.4. WERs and Latency on in-house Data

This section contains the results of the proposed model on large-scale in-house data. We use a 3-layer Conformer as the text encoder with a dropout rate of 0.1. Compared to the fast-slow RNN-T baseline, the deliberation model obtains around 3% WERR on ‘VA2’ while no improvement on ‘VA1’ where utterances have less than five words on average. The latter result is consistent with the analysis in the above section. Similar to Librispeech, on average, there is only a small extra emission delay introduced by the deliberation model. P95 emission delays are the same for all models we experiment with.

Table 4. WERs (%) and ED (ms) from baselines RNN-Ts and proposed fast-slow encoder based transducer on in-house data.

model	context	WER		ED	
		VA1	VA2	Avg	P95
RNN-T	160	4.73	12.89	390	560
fast-slow		4.65	12.43	369	560
+ deliberation	160/800	4.70	12.13	373	560
+ mask p = 0.1		4.68	12.08	373	560

5. CONCLUSION AND FUTURE WORK

This paper introduces a streaming deliberation model for fast-slow encoder based transducer, which improves the recognition accuracy of fast-slow encoder based transducer by leveraging partial hypotheses from the fast encoder. The proposed deliberation model shows 3-5% WERR on both Librispeech and in-house data compared to fast-slow encoder based transducer, and 10-11% WERR on Librispeech and up to 6% WERR on in-house data compared with baseline RNN-Ts. The proposed deliberation model introduces negligible extra emission delays. In the future, we will explore alternative text augmentation approaches for further accuracy improvement.

6. REFERENCES

- [1] Alex Graves, “Sequence transduction with recurrent neural networks,” *arXiv preprint arXiv:1211.3711*, 2012.
- [2] Kanishka Rao, Haşim Sak, and Rohit Prabhavalkar, “Exploring architectures, data and units for streaming end-to-end speech recognition with RNN-transducer,” in *Proc. of ASRU*, 2017.
- [3] Yanzhang He, Tara N Sainath, Rohit Prabhavalkar, Ian McGraw, Raziq Alvarez, Ding Zhao, David Rybach, Anjali Kannan, Yonghui Wu, Ruoming Pang, et al., “Streaming end-to-end speech recognition for mobile devices,” in *Proc. of ICASSP*, 2019.
- [4] Jay Mahadeokar, Yuan Shangguan, Duc Le, Gil Keren, Hang Su, Thong Le, Ching-Feng Yeh, Christian Fuegen, and Michael L Seltzer, “Alignment restricted streaming recurrent neural network transducer,” in *Proc. of SLT*, 2021.
- [5] Alex Graves and Navdeep Jaitly, “Towards end-to-end speech recognition with recurrent neural networks,” in *Proc. of ICML*, 2014.
- [6] Shinji Watanabe, Takaaki Hori, Suyoun Kim, John R Hershey, and Tomoki Hayashi, “Hybrid CTC/attention architecture for end-to-end speech recognition,” *IEEE Journal of Selected Topics in Signal Processing*, 2017.
- [7] Dario Amodei, Sundaram Ananthanarayanan, Rishita Anubhai, Jingliang Bai, Eric Battenberg, Carl Case, Jared Casper, Bryan Catanzaro, Qiang Cheng, Guoliang Chen, et al., “Deep speech 2: End-to-end speech recognition in English and Mandarin,” in *Proc. of ICML*, 2016.
- [8] William Chan, Navdeep Jaitly, Quoc Le, and Oriol Vinyals, “Listen, attend and spell: A neural network for large vocabulary conversational speech recognition,” in *Proc. of ICASSP*, 2016.
- [9] Jinyu Li, Yu Wu, Yashesh Gaur, Chengyi Wang, Rui Zhao, and Shujie Liu, “On the comparison of popular end-to-end models for large scale speech recognition,” in *Proc. of Interspeech*, 2020.
- [10] Tara N Sainath, Ruoming Pang, David Rybach, Yanzhang He, Rohit Prabhavalkar, Wei Li, Mirkó Visontai, Qiao Liang, Trevor Strohman, Yonghui Wu, et al., “Two-pass end-to-end speech recognition,” *arXiv preprint arXiv:1908.10992*, 2019.
- [11] Kevin Hu, Rohit Prabhavalkar, Ruoming Pang, and Tara Sainath, “Deliberation model based two-pass end-to-end speech recognition,” in *Proc. of ICASSP*, 2020.
- [12] Qiuqia Li, Chao Zhang, and Philip C Woodland, “Integrating source-channel and attention-based sequence-to-sequence models for speech recognition,” in *Proc. of ASRU*, 2019.
- [13] Jinxi Guo, Tara Sainath, and Ron Weiss, “A spelling correction model for end-to-end speech recognition,” in *Proc. of ICASSP*, 2019.
- [14] Arun Narayanan, Tara N Sainath, Ruoming Pang, Jiahui Yu, Chung-Cheng Chiu, Rohit Prabhavalkar, Ehsan Variiani, and Trevor Strohman, “Cascaded encoders for unifying streaming and non-streaming ASR,” in *Proc. of ICASSP*, 2021.
- [15] Bo Li, Anmol Gulati, Jiahui Yu, Tara N Sainath, Chung-Cheng Chiu, Arun Narayanan, Shuo-Yiin Chang, Ruoming Pang, Yanzhang He, James Qin, et al., “A better and faster end-to-end model for streaming ASR,” in *Proc. of ICASSP*, 2021.
- [16] Jay Mahadeokar, Yangyang Shi, Ke Li, Duc Le, Jiedan Zhu, Vikas Chandra, Ozlem Kalinli, and Michael L. Seltzer, “Streaming parallel transducer beam search with fast-slow cascaded encoders,” in *Proc. of Interspeech*, 2022.
- [17] Ke Hu, Tara N Sainath, Arun Narayanan, Ruoming Pang, and Trevor Strohman, “Transducer-based streaming deliberation for cascaded encoders,” in *Proc. of ICASSP*, 2022.
- [18] Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, and Ruoming Pang, “Conformer: Convolution-augmented transformer for speech recognition,” in *Proc. of Interspeech*, 2020.
- [19] Ke Hu, Tara N. Sainath, Yanzhang He, Rohit Prabhavalkar, Trevor Strohman, Sepand Mavandadi, and Weiran Wang, “Improving deliberation by text-only and semi-supervised training,” in *Proc. of Interspeech*, 2022.
- [20] Weiran Wang, Ke Hu, and Tara N Sainath, “Deliberation of streaming RNN-transducer by non-autoregressive decoding,” in *Proc. of ICASSP*, 2022.
- [21] Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur, “Librispeech: an ASR corpus based on public domain audio books,” in *Proc. of ICASSP*, 2015.
- [22] Daniel S Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D Cubuk, and Quoc V Le, “SpecAugment: A simple data augmentation method for automatic speech recognition,” *arXiv preprint arXiv:1904.08779*, 2019.
- [23] Yangyang Shi, Yongqiang Wang, Chunyang Wu, Ching-Feng Yeh, Julian Chan, Frank Zhang, Duc Le, and Mike Seltzer, “Emformer: Efficient memory transformer based acoustic model for low latency streaming speech recognition,” in *Proc. of ICASSP*, 2021.