

LEGO-FEATURES: EXPORTING MODULAR ENCODER FEATURES FOR STREAMING AND DELIBERATION ASR

Rami Botros, Rohit Prabhavalkar, Johan Schalkwyk, Ciprian Chelba, Tara N. Sainath, Françoise Beaufays

Google LLC, USA

ramibotros@google.com

ABSTRACT

In end-to-end (E2E) speech recognition models, a representational tight-coupling inevitably emerges between the encoder and the decoder. We build upon recent work that has begun to explore building encoders with modular encoded representations, such that encoders and decoders from different models can be stitched together in a zero-shot manner without further fine-tuning. While previous research only addresses full-context speech models, we explore the problem in a streaming setting as well. Our framework builds on top of existing encoded representations, converting them to modular features, dubbed as *Lego-Features*, without modifying the pre-trained model. The features remain interchangeable when the model is re-trained with distinct initializations. Though sparse, we show that the Lego-Features are powerful when tested with RNN-T or LAS decoders, maintaining high-quality downstream performance. They are also rich enough to represent the first-pass prediction during two-pass deliberation. In this scenario, they outperform the N-best hypotheses, since they do not need to be supplemented with acoustic features to deliver the best results. Moreover, generating the Lego-Features does not require beam search or auto-regressive computation. Overall, they present a modular, powerful and cheap alternative to the standard encoder output, as well as the N-best hypotheses.

Index Terms— modular, representations, zero-shot stitching

1. INTRODUCTION

E2E speech recognition models, which combine acoustic, pronunciation and language models from conventional systems [1] into one neural network, have become widely used, especially for on-device applications [2, 3, 4, 5, 6, 7]. Since they are much smaller than conventional models, and their inference speed is often much faster [2, 3, 8, 9], they work well for various streaming applications. They typically use an encoder-decoder architecture [10]. Like most deep neural networks, the whole architecture is usually trained end to end. The encoder implicitly learns to serve the subsequent decoder layers, and thus conversely, the decoder is thoroughly oriented towards inputs coming from the specific encoder that it has been trained with. Therefore, encoders and decoders from different models or training runs, are generally not interchangeable without further E2E training.

This tight coupling between both components stands in the way of a flexible, modular architecture. Speech encoders that have been trained on high-resource ASR data can serve as foundation models for other tasks like sentiment analysis [11] or low-resource translation [12], to name a few. However, this presents a challenge if a shared encoder representation is used for multiple downstream tasks: When the ASR encoder is retrained, all downstream models must be retrained as well. Hence, it would be more practical if each component can be developed and updated independently. To that end,

we present a method for building modular speech encoder features, where different versions of the encoder can be plugged into the decoder in a zero-shot stitching manner without fine-tuning.

Our method works by building on top of an existing base encoder, which is kept frozen. We adapt the Beam-Convolution scheme described in [13] to train streaming modular encoded representations, which we call Lego-Features. To produce them, the original (fixed) continuous encoded features pass through a few extra trainable “Exporter” layers, then through a CTC decoder, which is trained with an auxiliary CTC loss. Lego-Features are defined as the sorted top K CTC logit indices at every frame, see Figure 1. The logits operate over a discrete space (here: wordpiece vocabulary) and are grounded in the transcript text, which is why they tend to be modular. Overall, the traditional encoder features are forced through a tight discretizing bottleneck, which protects downstream models from coupling themselves to fine details in the encoded representation. Downstream consumers of Lego-Features need to first re-embed them, since they come in as sparse indices.

[13, 14] have shown how this tight bottleneck still produces a powerful representation which is sufficiently informative for downstream ASR decoders. They also perform a “modularity test”: The downstream decoder is kept constant, but gets input with a new version of the encoded representation, which is obtained by retraining the encoder from scratch using a different initialization. The switch is done in a zero-shot manner without any extra fine-tuning. Traditional continuous encoded features categorically fail the modularity test, bringing the downstream performance to nearly 100% WER, which is what motivates this new type of encoded representation. We build on the original works with a few novel contributions:

- 1) We find that training the modular encoder from scratch under the CTC loss is insufficient for producing the best performance. Instead, our recipe pre-trains some base encoder layers with RNN-T loss and keeps them frozen. Next, we just train the extra Exporter layers with the auxiliary CTC loss. This solution is also practical since it enables researchers to cheaply export modular features without having to modify their original system. Thus, the quality, latency and efficiency of the base model are all maintained.

- 2) We adapt the design to a streaming setting for the first time. Unlike the original work [13, 14], our encoder layers attention have limited left and right context windows, and the produced Lego-Features are successfully paired with a streaming-friendly RNN-T decoder. The streaming architecture still exhibits strong downstream ASR quality and passes the modularity test. By plugging the same fixed set of Lego-Features into causal as well as non-causal decoders, our work adds further evidence to their modularity and interoperability.

- 3) Rather than merely looking at the Lego-Features as an encoded representation, we also study them as an alternative to the N-best hypotheses within two-pass systems. We provide new comparisons against the N-best in terms of speed, accuracy and modularity. To

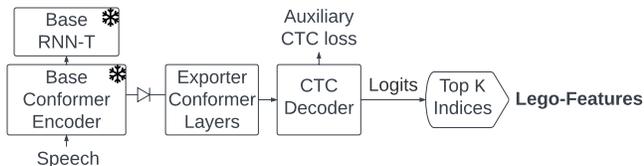


Fig. 1. Modular Encoder. Lego-Features are exported from frozen base encoder by training extra layers with an auxiliary CTC loss.

this end, the Lego-Features are used as a first-pass output within the deliberation framework [15]. This achieves good post-deliberation WER performance, which is shown to be on-par with a baseline that performs deliberation on 1st-pass RNN-T N-best hypotheses + audio features. The Lego-Features demonstrate success in the modularity test here as well. On the other hand, we find that the N-best hypothesis text does not pass the modularity test, i.e. a new N-best from a second model would confuse the deliberation decoder from the first, which is a novel observation. Moreover, the Lego-Features are cheaper to produce than the N-best, since they require no beam-search or auto-regressive decoding, but are generated via a simple projection at every frame.

Other works have attempted to present generic methods for zero-shot stitching between layers. In [16], this is achieved by learning representations relative to data-dependent anchors. In contrast, the method presented here does not need to choose anchor samples and leverages the existence of ground-truth speech transcripts instead. Another general approach, presented in [17], uses self-supervised objectives designed to encourage compatibility of different layer outputs. It is an open question whether the cited methods can deal with long sequences, whereas the CTC loss used here is a natural choice that works well with ASR and gives interpretable outputs.

Further, some research has already experimented with deliberation on top of CTC outputs to save the cost of first-pass decoding [18, 19, 20]. This includes the Align-refine approach, which iteratively improves on the first-pass output. Those works tend to focus on optimizing the size and speed of the first-pass model, whereas our focus is mainly on modularity. Nevertheless, since we build on base encoder layers that have been pre-trained with the RNN-T loss, we find our CTC outputs to have high quality, which removes the need for audio attention that is used in other deliberation models. Hence, this work also introduces some speed gains to deliberation, without using the iterative Align-refine approach.

On the whole, with one simple representations, we get a compelling cheap, streaming-friendly, as well as modular, alternative to both the continuous encoding vector and the N-best hypotheses, without any loss in quality.

2. MODELING

Our framework is trained in three separate stages described below.

2.1. Base Model

We start off from a pre-trained end-to-end system that follows the cascade architecture in [21]: The base encoder comprises 3 convolution layers, then 14 Conformer [22] blocks: 4 causal ones, followed by 5 blocks that process 180 milliseconds of right-context each, then 5 more causal ones. This base encoder is pre-trained using the RNN-T loss on the same training set. For the modularization steps below,

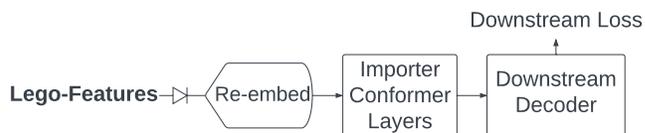


Fig. 2. Downstream models embed and process the fixed Lego-features before passing them to a downstream decoder.

the pre-trained RNN-T decoder layers will be discarded, and the base encoder is kept frozen. This recipe allows us to keep the existing pre-trained model unchanged while exporting modular features.

2.2. Exporting Lego-Features

Figure 1 shows how the modular encoder is trained on top of a frozen base model. The Exporter layers comprise further Conformer blocks with 180ms look-ahead context. The CTC decoder [23] amounts to a single projection layer to compute the frame-level posterior over the output vocabulary. Our work uses wordpiece output tokens, but further research can explore using phonemes or graphemes instead. The depicted CTC loss is applied to those logits and is what trains the Exporter layers. Finally, the Lego-Features are computed by extracting the sorted top- K indices of the CTC logits, giving K integers at every frame. Note that this is performed on the logit vector directly, without requiring any actual decoding algorithm like beam-search.

2.3. Downstream Models

Figure 2 illustrates how downstream models generally consume the Lego-Features, which come in as sparse indices. The downstream consumer does not receive extra information about how the indices map to wordpiece tokens, and hence starts by embedding them. An Importer module, once again consisting of 180ms look-ahead Conformer blocks, prepares the embeddings for the downstream decoder. [13, 14] use 1D convolution + multi-headed attention in place of the Importer, but our early experiments show that Conformer blocks improve over this original stack. Note that the Lego-Features themselves are kept constant during downstream training. We experiment with two types of ASR decoders as examples for downstream tasks, which are used with the same fixed set of Lego-Features.

2.3.1. Downstream RNN-T Decoder

The first downstream model uses an RNN-T decoder, which tends to serve real-time applications well, since it processes the input frames in a streaming fashion as they become available and starts outputting text tokens after a short delay [3, 24]. We adopt the same RNN-T decoder layer architecture from the base model (Section 2.1) but use it as a simulated downstream task, as the decoder in Figure 2, to see if the bottlenecked Lego-Features are as informative as the continuous base encoded tensor.

2.3.2. Downstream LAS decoder / Deliberation

As a second downstream ASR decoder in Figure 2, we experiment with a full-context Listen-Attend-and-Spell (LAS) decoder [25], which can achieve higher quality by attending to all input frames.

A fitting baseline to this experiment is second-pass deliberation ASR [15]. Typically, a deliberation system generates first-pass

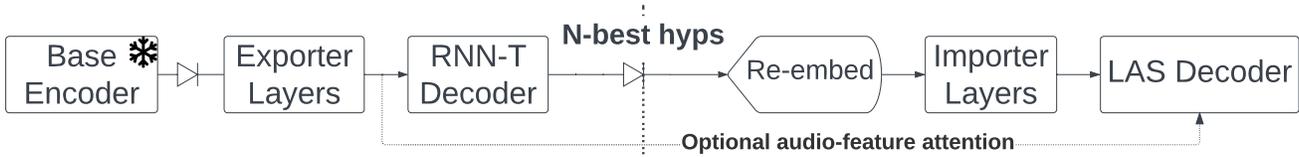


Fig. 3. Baseline deliberation on N-best RNN-T hyps. The LAS decoder attends to embedded text and optionally to the pre-RNN-T audio features. Modularity test boundary shown as the dotted line in the middle.

hypotheses using a fast decoder, like RNN-T, then embeds its N-best hyps and attends to them with a second-pass full-context LAS decoder. We have therefore constructed a comparable deliberation baseline model shown in Figure 3. This model is analogous to our full pipeline, i.e. Figures 1 & 2 put together, and is designed to have a similar total model size and encoder latency. It starts with the same frozen base encoder, then trains a first-pass RNN-T decoder to obtain the N-best hyps, which stands to be compared to the Lego-Features in terms of informativeness and modularity. Figure 3 also ends with an LAS decoder, except this one can optionally attend to the continuous encoder features as well, as is done in previous deliberation work [15]. Gradients do not flow back through embedded N-best.

3. EXPERIMENTAL SETTINGS

3.1. CTC Logit Evaluation

An interesting aspect of the Lego-Features encoder is that one can evaluate its quality directly before providing the features to any downstream tasks. This is done via a preliminary experiment where we directly decode from the full set of the CTC-trained logits (before the top- K operation in Figure 1) using beam search or greedy decoding. The decoding algorithm used for this evaluation is tangential to how the Lego-Features are produced, since those are only extracted as the top- K logit ranks without decoding actual transcripts. Yet this direct evaluation can inform us about the general quality of the CTC-trained logits, from which the Lego-Features are produced.

3.2. WER and Modularity Test

The downstream ASR decoders trained on the Lego-Features (Section 2.3) are then evaluated and a modularity test is performed. The aim of the test is to check if two different versions of the encoded features are interchangeable. We test that by keeping the downstream model fixed, but feeding it with a new version of the encoded features, which we get from another training run. The second training is done from scratch with a new initialization. We compare the WER performance of the decode before and after the switch, denoted as “Normal \rightarrow Mod. Test WER” in our tables. For the Lego-Features, we retrain the encoder in Figure 1, where the base frozen encoder is also replaced with a second version from a retrained base. As a baseline, we also test the modularity of the base model itself, where we simply train the base encoder + decoder a second time end-to-end and get the retrained encoder from there.

3.3. Architectural Details

Our base architecture follows [21]: All Conformer layers [22] are 512-dim, use 8-headed self-attention and a convolution kernel size

of 15. We train on a 128D log-mel feature frontend with a 16-D one-hot domain-id vector appended to it, see [26].

Our models work with 4,096 word pieces [27]. The RNN-T decoder comprises a prediction network and a joint network with a single 640-dim FF layer. The embedding prediction network [28], uses an embedding dimension of 320, and has 9M parameters. For the deliberation decoder, we use a 2-layer LSTM similar to [15], where each layer has 1536 hidden units followed by 384-dim projection. We do not use external LMs.

3.4. Datasets

As discussed in [29], all E2E models are trained on multidomain audio-text pairs [26]. All datasets obtain their labels in a semi-supervised fashion, using larger teacher models trained on in-domain data to provide pseudo labels [30, 31]. Data was handled in accordance to Google AI principles [32]. To further increase data diversity, multi-condition training (MTR) [33], random data down-sampling to 8kHz [34] and SpecAug [35] are also used. Noisy data is generated at signal-noise-ratio (SNR) from 0 to 30 dB, with an average SNR of 12 dB, and with T60 times ranging from 0 to 900ms, averaging 500ms. Noise segments are sampled from YouTube and daily life noisy environmental recordings. Both 8 kHz and 16 kHz versions of the data are generated, each with equal probability, to make the model robust to varying sample rates.

The *Voice-Search* test set has 10K Voice Search utterances with an average length of 5.5 seconds. They are anonymized, hand-transcribed, and are representative of Google’s Voice Search traffic.

4. EXPERIMENTAL RESULTS

4.1. Preliminary CTC Decoder Evaluation

Exporter Properties			CTC Test WER	
# Blocks	Size	Right Context	Greedy	Beam-search (Oracle)
1	10M	+180 ms	5.9%	5.8% (2.8%)
3	30M	+540 ms	5.5%	5.3% (2.7%)

Table 1. CTC Voice-Search WER for different Exporter setups

As explained in Section 3.1, the CTC decoder in Figure 1 can be evaluated directly. Table 1 shows two settings for the Exporter layers and their corresponding CTC WER performance. The right-context length indicates the extra duration of future context attended to by the Exporter, noting that the base encoder already sees a future context of 900ms. In both cases, greedy decoding performs close

Type	Encoder		RNN-T WER
	Size	Right-Context	Normal → Mod. Test
Base	146M	900 ms	6.4% → 99%
Modularized	207M	1440 ms	5.6% → 5.6%

Table 2. Downstream RNN-T Test WER with Modularity Test. The base encoder is from the original pre-trained model.

to beam search, which tracks 16 hypotheses in its beam. For all the downstream experiments below, we use the better setup with 3 blocks for the Exporter, and apply the same design to the Importer.

4.2. Base RNN-T vs. Downstream RNN-T

Our first downstream setting works with an RNN-T decoder (Section 2.3.1). Table 2 demonstrates how the Lego-Features bottleneck still produces a rich encoding that the downstream Importer and RNN-T use well. We export $K = 12$ Lego-Features per frame and the downstream re-embeds each into 32 dimensions. Preliminary experiments, omitted here for brevity, indicate that varying these values does not affect downstream WER performance significantly. The Base case in the table is simply the frozen base model on the left of Figure 1, in which case the modularity test connects a new base encoder (from another training run) to the same frozen base RNN-T. The modularity test fails for the base case, yet passes for the Lego-Features. Both models involve different sizes and latencies, so a direct WER contest between them is not the main concern. Rather, the goal is to show that the Lego-Features bottleneck does not degrade performance while enabling modularity.

To test robustness across changing domains, we also supply the same Lego-Features used above to a downstream RNN-T model that is trained on Librispeech data instead. The modularity test results are shown in Table 3 and only cause less than 4% relative WER decline.

4.3. Deliberation on N-Best vs. Lego-Features

Table 4 compares the LAS deliberation scenarios described in Section 2.3.2, where the Lego-Features are compared to an N-best as a first-pass output. Dropping the audio connection significantly degrades performance in the N-best case, which is consistent with previous findings [15]. The Lego-Features seem to preserve more information in the encoding, and thus do not need the audio connection. They are significantly better than N-best text, and are only off by 0.1 in absolute WER from N-best + audio.

The modularity test causes no performance decline for the Lego-Features, but does not work well in the N-best case; even the text-only case degrades by 17% relative WER. This somewhat unexpected result might be a symptom of label bias, which RNN-T suffers from because of local normalization [36, 37], but the CTC decoder avoids with its conditional independence assumption. Hence, two separately-trained RNN-T first-pass models might exhibit different biases in their N-bests, leading to this result.

Dev-Clean WER	Test-Other WER
Normal → Mod. Test	Normal → Mod. Test
4.9 → 5.1	10.0 → 10.3

Table 3. Modularity tests if downstream is trained on Librispeech

First Pass	Embedded Shape	Attend	Downstream WER
		Audio	Normal → Mod. Test
RNN-T N -best	$[N \cdot U, E_1]$	No Yes	5.4% → 6.3% 5.0% → 14.3%
Lego-Features	$[T, K \cdot E_2]$	No	5.1% → 5.1%

Table 4. Deliberation WER and Modularity Tests. Embedded Shapes discussed in Section 4.3.1

4.3.1. Speed Comparison

Table 4 notes a difference in the input shapes to the Importers across the different types of first-pass models, after re-embedding in Figure 2 & 3. Here, E_1 and E_2 are the respective embedding dimensions, n is the RNN-T’s beam width and U is the number of text tokens produced by it. K is the number of logit indices in the Lego-Features and T is their sequence length (=number of encoded frames). Note how the N-best’s embedding expands the output sequence length, since it stacks the N hypothesis sequentially while keeping the sentence structures intact, in order to attend to this order information during second-pass decoding. Since the LegoFeatures describe per-frame logit ranks without serializing them into sentences, we forgo this expansion and concatenate the embeddings within the depth dimension at each frame instead. This saves on computational cost, since the #GFLOPs used by LAS is proportional to the sequence length it is attending to. While U can change from one utterance to the other, the embedded matrices have to be padded to maximum length when working with hardware accelerators. Our system uses $n = 8$, $U = 120$, $E_1 = 384$, $T = 343$, $K = 12$, and $E_2 = 32$. This makes the depth dimension equal, but LegoFeatures’ sequence length is 64% smaller than the N-best’s.

Another important computational benefit of deliberating on LegoFeatures is that we can obtain them without performing a beam-search procedure. It is hence possible to compute them for long utterances with high parallelization, only limited by the number of TPU cores available. Generating the N-best, on the other hand, requires sequential auto-regressive processing. For instance, benchmarking this sequential path in the RNN-T (using an in-house server TPU and the above dimensions) gives 1.8 ms per output token, or 216 ms per utterance in the padded worst case, which does become the bottleneck after the other layers are parallelized.

5. CONCLUSIONS AND FUTURE WORK

In this paper, we describe a simple recipe for exporting streaming-friendly modular encoded representations and successfully test them with RNN-T and LAS decoders. Overall, exporting the encoder output as top CTC-trained logits introduces multiple benefits. The encoding achieves strong WER performance and interchangability is demonstrated through the modularity test. If regarded as a representation for first-pass ASR prediction, the Lego-Features surpass the N-best in quality, modularity, and generation speed.

To address resource-limited environments like on-device ASR, and to improve latency, future research can explore using smaller Exporter and Importer layers. Another avenue is to export CTC logits over phoneme/triphone/grapheme vocabularies, or a combination thereof. Different types of Lego-Features can be tested with various downstream tasks, like confidence models, speech translation or spoken language understanding.

6. REFERENCES

- [1] G. Pundak and T. N. Sainath, “Lower frame rate neural network acoustic models,” in *Proc. Interspeech*, 2016.
- [2] B. Li, A. Gulati, J. Yu, et al., “A Better and Faster End-to-End Model for Streaming ASR,” in *Proc. ICASSP*, 2021.
- [3] Y. He, T. N. Sainath, R. Prabhavalkar, et al., “Streaming End-to-end Speech Recognition For Mobile Devices,” in *Proc. ICASSP*, 2019.
- [4] C.-C. Chiu, T. N. Sainath, Y. Wu, et al., “State-of-the-art Speech Recognition With Sequence-to-Sequence Models,” in *Proc. ICASSP*, 2018.
- [5] S. Kim, T. Hori, and S. Watanabe, “Joint CTC-attention based end-to-end speech recognition using multi-task learning,” in *Proc. ICASSP*, 2017.
- [6] J. Li, R. Zhao, H. Hu, and Y. Gong, “Improving RNN transducer modeling for end-to-end speech recognition,” in *Proc. ASRU*, 2019.
- [7] A. Zeyer, A. Merboldt, R. Schlüter, and H. Ney, “A new training pipeline for an improved neural transducer,” in *Proc. Interspeech*, 2020.
- [8] T. N. Sainath, Y. He, Narayanan, et al., “An Efficient Streaming Non-Recurrent On-Device End-to-End Model with Improvements to Rare-Word Modeling,” in *Interspeech*, 2021.
- [9] Tara N Sainath, Yanzhang He, Bo Li, Arun Narayanan, et al., “A streaming on-device end-to-end model surpassing server-side conventional model quality and latency,” in *Proc. ICASSP*. IEEE, 2020, pp. 6059–6063.
- [10] R. Prabhavalkar, K. Rao, T. N. Sainath, B. Li, L. Johnson, and N. Jaitly, “A comparison of sequence-to-sequence models for speech recognition,” in *Proc. Interspeech*, 2017.
- [11] Zhiyun Lu, Liangliang Cao, Yu Zhang, Chung-Cheng Chiu, and James Fan, “Speech sentiment analysis via pre-trained features from end-to-end asr models,” in *Proc. ICASSP*. IEEE, 2020, pp. 7149–7153.
- [12] Sameer Bansal, Herman Kamper, Karen Livescu, Adam Lopez, and Sharon Goldwater, “Pre-training on high-resource speech recognition improves low-resource speech-to-text translation,” *arXiv preprint arXiv:1809.01431*, 2018.
- [13] Siddharth Dalmia, Abdelrahman Mohamed, Mike Lewis, Florian Metze, and Luke Zettlemoyer, “Enforcing encoder-decoder modularity in sequence-to-sequence models,” *arXiv preprint arXiv:1911.03782*, 2019.
- [14] Siddharth Dalmia, Dmytro Okhonko, Mike Lewis, Sergey Edunov, Shinji Watanabe, Florian Metze, Luke Zettlemoyer, and Abdelrahman Mohamed, “Legonn: Building modular encoder-decoder models,” *arXiv:2206.03318*, 2022.
- [15] Ke Hu, Tara N Sainath, Ruoming Pang, and Rohit Prabhavalkar, “Deliberation model based two-pass end-to-end speech recognition,” in *ICASSP*. IEEE, 2020, pp. 7799–7803.
- [16] Luca Moschella, Valentino Maiorca, Marco Fumero, Antonio Norelli, Francesco Locatello, and Emanuele Rodolà, “Relative representations enable zero-shot latent space communication,” *arXiv preprint arXiv:2209.15430*, 2022.
- [17] Michael Gygli, Jasper Uijlings, and Vittorio Ferrari, “Towards reusable network components by learning compatible representations,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2021, vol. 35, pp. 7620–7629.
- [18] Ethan A Chi, Julian Salazar, and Katrin Kirchhoff, “Align-refine: Non-autoregressive speech recognition via iterative re-alignment,” *arXiv preprint arXiv:2010.14233*, 2020.
- [19] Weiran Wang, Ke Hu, and Tara N Sainath, “Deliberation of streaming rnn-transducer by non-autoregressive decoding,” in *Proc. ICASSP*. IEEE, 2022, pp. 7452–7456.
- [20] Weiran Wang, Ke Hu, and Tara N Sainath, “Streaming align-refine for non-autoregressive deliberation,” *arXiv preprint arXiv:2204.07556*, 2022.
- [21] Tara N Sainath, Yanzhang He, Arun Narayanan, Rami Botros, et al., “Improving the latency and quality of cascaded encoders,” in *Proc. ICASSP*. IEEE, 2022, pp. 8112–8116.
- [22] A. Gulati, J. Qin, C.-C. Chiu, et al., “Conformer: Convolution-augmented Transformer for Speech Recognition,” in *Proc. Interspeech*, 2020.
- [23] A. Graves, S. Fernandez, F. Gomez, and J. Schmidhuber, “Connectionist Temporal Classification: Labeling Unsegmented Sequence Data with Recurrent Neural Networks,” in *Proc. ICML*, 2006.
- [24] A. Graves, “Sequence Transduction with Recurrent Neural Networks,” *CoRR*, vol. abs/1211.3711, 2012.
- [25] W. Chan, N. Jaitly, Q. V. Le, and O. Vinyals, “Listen, attend and spell,” *CoRR*, vol. abs/1508.01211, 2015.
- [26] A. Narayanan, R. Prabhavalkar, C.-C. Chiu, et al., “Recognizing Long-Form Speech Using Streaming End-to-End Models,” in *Proc. ASRU*, 2019.
- [27] M. Schuster and K. Nakajima, “Japanese and Korean voice search,” in *Proc. ICASSP*, 2012.
- [28] R. Botros and T.N. Sainath, “Tied & reduced rnn-t decoder,” in *Proc. Interspeech*, 2021.
- [29] T. N. Sainath, Y. He, B. Li, et al., “A Streaming On-Device End-To-End Model Surpassing Server-Side Conventional Model Quality and Latency,” in *Proc. ICASSP*, 2020.
- [30] D. Hwang, K. Sim, Z. Huo, and T. Strohmaier, “Pseudo Label Is Better Than Human Label,” in *Proc. ICASSP*, 2022.
- [31] H. Liao, E. McDermott, and A. Senior, “Large Scale Deep Neural Network Acoustic Modeling with Semi-supervised Training Data for YouTube Video Transcription,” in *Proc. ASRU*, 2013.
- [32] “Google ai principles,” <https://ai.google/principles/>.
- [33] C. Kim, A. Misra, K. Chin, et al., “Generation of Large-Scale Simulated Utterances in Virtual Rooms to Train Deep-Neural Networks for Far-Field Speech Recognition in Google Home,” in *Proc. Interspeech*, 2017.
- [34] J. Li, D. Yu, J. Huang, and Y. Gong, “Improving Wideband Speech Recognition using Mixed-bandwidth Training Data in CD-DNN-HMM,” in *Proc. SLT*, 2012.
- [35] D. S. Park, W. Chan, Y. Zhang, C. Chiu, B. Zoph, E.D. Cubuk, and Q.V. Le, “SpecAugment: A Simple Data Augmentation Method for Automatic Speech Recognition,” in *Proc. Interspeech*, 2019.
- [36] Awni Hannun, “The label bias problem,” 2020.
- [37] Brian Yan, Siddharth Dalmia, Yosuke Higuchi, Graham Neubig, Florian Metze, Alan W Black, and Shinji Watanabe, “Ctc alignments improve autoregressive translation,” *arXiv preprint arXiv:2210.05200*, 2022.