

A Multi-level Acoustic Feature Extraction Framework for Transformer Based End-to-End Speech Recognition

Jin Li^{1,2}, Rongfeng Su^{1,2}, Xurong Xie^{1,2,3}, Nan Yan^{1,2}, Lan Wang^{1,2}

¹ CAS Key Laboratory of Human-Machine Intelligence-Synergy Systems,

² Guangdong-Hong Kong-Macao Joint Laboratory of Human-Machine Intelligence-Synergy Systems, Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences, Shenzhen, China,

³ Institute of Software, Chinese Academy of Sciences, Beijing, China

{li.jin, rf.su, nan.yan, lan.wang}@siat.ac.cn, xurong@iscas.ac.cn

Abstract

Transformer based end-to-end modelling approaches with multiple stream inputs have been achieved great success in various automatic speech recognition (ASR) tasks. An important issue associated with such approaches is that the intermediate features derived from each stream might have similar representations and thus it is lacking of feature diversity, such as the descriptions related to speaker characteristics. To address this issue, this paper proposed a novel multi-level acoustic feature extraction framework that can be easily combined with Transformer based ASR models. The framework consists of two input streams: a shallow stream with high-resolution spectrograms and a deep stream with low-resolution spectrograms. The shallow stream is used to acquire traditional shallow features that is beneficial for the classification of phones or words while the deep stream is used to obtain utterance-level speaker-invariant deep features for improving the feature diversity. A feature correlation based fusion strategy is used to aggregate both features across the frequency and time domains and then fed into the Transformer encoder-decoder module. By using the proposed multi-level acoustic feature extraction framework, state-of-the-art word error rate of 21.7% and 2.5% were obtained on the HKUST Mandarin telephone and Librispeech speech recognition tasks respectively.

Index Terms: end-to-end speech recognition, transformer, feature extraction, multi-stream

1. Introduction

In recent years, end-to-end based modelling approaches have gained popularity in automatic speech recognition (ASR) community [1, 2, 3, 4, 5, 6]. End-to-end ASR models simplify traditional hybrid ASR models by using one single deep neural network instead of acoustic, pronunciation and language components, and thus text can be transcribed directly from speech. Previous researches have shown that significant system performance can be obtained from end-to-end ASR models over hybrid ASR models [7, 8].

Transformer [9] with self-attention mechanism is one of the effective end-to-end ASR architectures. The self-attention mechanism learns the temporal contextual information of the input sequence by applying attention matrices on the input frames. The motivation behind the self-attention mechanism in Transformer is to aggregate long-term dependencies among acoustic features in the encoder and compute the output sequences in parallel. Current Transformer based ASR modelling approaches can be classified into two categories. The first category is single-stream based methods [2, 4, 6, 10, 11, 12]. Most

studies of this category focus on exploring new neural network architecture for faster training speed and online decoding, but they ignore the different characteristic representations of the input signal. In contrast, the multi-stream based modelling approaches mainly focus on the front-end feature processing by using appropriate network structures. For example, in order to enrich the feature diversity, Wang et al. [13] designed a hierarchical attention mechanism to dynamically combine the knowledge from parallel streams. More recently, Chang et al. [14] used three different attention layers to acquire spectral and spatial feature representations from multiple source channels. However, since the previous works of multi-stream based approaches use the same resolution spectrograms for each stream as inputs, the intermediate features derived from each stream might have similar representations and lack diverse feature expressions, such as the vocal characteristics for determining different speakers.

To address this issue, this paper proposes a novel multi-level acoustic feature extraction framework for Transformer based ASR models. Inspired by [15], this framework consists of two input streams: a shallow stream with high-resolution spectrograms and a deep stream with low-resolution spectrograms. The shallow branch with less network layers is suggested to preserve the detailed information in the audio inputs that is beneficial for the classification of phones or words. The deep branch with more network layers is suggested to capture the utterance-level speaker characteristic information to improve the feature diversity. Inspired by the non-local operations from [16], a feature correlation based fusion strategy is proposed to integrate both features across the frequency and time domains. The fused features are then fed into Transformer encoder-decoder module to improve the ASR system performance.

The main contributions are summarized as below.

1. A novel multi-level acoustic feature extraction framework is proposed in this paper. Speaker-related information can be derived from the deep stream to improve the diversity of the feature representations of speech signal.
2. An additional feature correlation based fusion strategy is proposed to integrate two different types of feature expressions across the frequency and time domains.
3. The Transformer based ASR system using the proposed multi-level acoustic feature extraction framework gave state-of-the-art error rates of 21.7% and 2.5% on the HKUST Mandarin telephone speech [17] and Librispeech [18] recognition tasks respectively.

The rest of this paper is organized as follows. The methodology about the multi-level acoustic feature extraction frame-

work is proposed in Section 2. Experiments and results are presented in Section 3. Finally, conclusions and future works are drawn in Section 4.

2. Methodology

2.1. Transformer based ASR model

As illustrated in Figure 1, the Transformer based ASR model used in this paper consists of the proposed multi-level acoustic feature extraction module, as well as the Transformer encoder and decoder. The multi-level acoustic feature extraction module has two input streams, one is a deep stream with low-resolution spectrograms and the other is a shallow stream with high-resolution spectrograms. After getting features extracted from two streams, they will be integrated by using the fusion strategy proposed in Section 2.3 and then used the inputs of Transformer. Standard Transformer [9] containing encoder and decoder components is used in this paper. The encoder is composed of stacked layers that have identical components. Each layer in the encoder has two sub-layers. The first one is a multi-head self-attention sub-layer and the other is a simple, position wise fully connected feed-forward network with ReLU activation function. The positional encodings are suggested to use to capture longer dependencies in a sentence. The decoder structure is similar to the encoder. The only difference is: the decoder inserts a third sub-layer that performs multi-head attention over the output of the encoder stack. In addition, to enhance the information representations, residual connections around each of the sub-layers are used, which are followed by layer normalization.

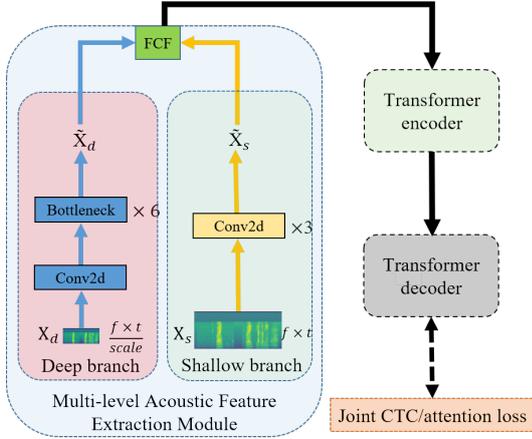


Figure 1: Overview of the Transformer based ASR model: “FCF” represents the fusion strategy presented in Figure 2.

2.2. Multi-level acoustic feature extraction

There are two streams with different spectrogram resolutions in the multi-level acoustic feature extraction module. Assumed that the input of the shallow stream is $\tilde{X}_s := (x_1, \dots, x_T) = \{x_t \in \mathbb{R}^D | t = 1, 2, \dots, T\}$ while the input of the deep stream is $\tilde{X}_d := (x'_1, \dots, x'_{T'}) = \{x'_t \in \mathbb{R}^{D'} | t = 1, 2, \dots, T'\}$, where T' and T ($T' < T$) represent the different scales along the time axis, D and D' are the low and high resolution size of spectrograms respectively. The high resolution spectrogram is suggested to provide the detailed acoustic information especially for the alignment between speech features and labels, which is

Table 1: Neural network configuration in the multi-level acoustic feature extraction module.

Stream	Block	Output Size	Repeats	Stride
Deep	Conv2d	32	1	2
	Bottleneck	32	1	1
	Bottleneck	32	1	1
	Bottleneck	48	3	2
	Bottleneck	64	3	2
	Bottleneck	128	2	1
Shallow	Conv2d	128	1	2
	Conv2d	256	1	2
	Conv2d	256	1	1

useful for classifying different phones or words. The low resolution spectrogram is suggested to contain the most important information about the speaker characteristics in each utterance.

The configuration of the neural network for processing the deep and shallow streams is shown in Table 1. The neural network for the deep stream contains one 2D convolution block followed by six bottleneck residual blocks from MobileNetV2 [19]. The neural network for the shallow stream consists of three 2D convolution blocks. Each convolution block involves batch normalization layer and ReLU activation function.

2.3. Fusion strategy

As mentioned before, the intermediate features derived from the shallow branch contain the detailed information that is helpful for distinguishing different phones or words, while the outputs of the deep branch contain the information for identifying different speakers. In order to remove the redundant parts between the shallow and deep features while preserve the complementary ones, a novel feature correlation based fusion strategy shown in Figure 2 is used to integrate the both intermediate features across the frequency and time domains.

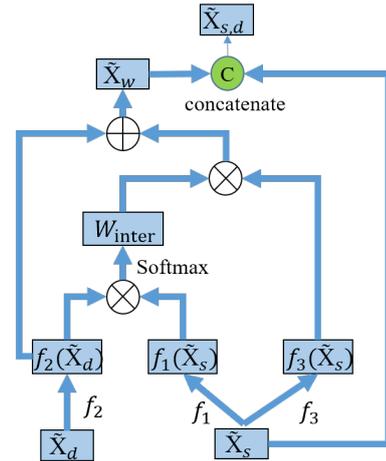


Figure 2: Feature Correlation based Fusion (FCF) strategy: The interaction weight W_{inter} is used to determine the correlation between the shallow feature \tilde{X}_s and deep feature \tilde{X}_d . “ \otimes ” represents the dot product operation.

Assumed $\tilde{X}_s \in \mathbb{R}^{c,d,m}$ represents the shallow stream feature vector, $\tilde{X}_d \in \mathbb{R}^{c,d,m}$ represents the deep stream feature

vector rescaled by bilinear interpolation of \tilde{X}_s , where c is the number of channels, $d_m = t \times f$ and f is frequency. The interaction weight between the both features is computed as follow.

$$W_{\text{inter}} := \text{Softmax}\left(\frac{f_1(\tilde{X}_s)f_2(\tilde{X}_d)^T}{\sqrt{d_m}}\right) \quad (1)$$

where $f_1(\cdot)$ and $f_2(\cdot)$ are two different linear transformations with 1D convolution operations, both the outputs of $f_1(\cdot)$ and $f_2(\cdot)$ are two-dimensional matrices, the ‘‘Softmax’’ activation function across rows is used. In our early experiments, a weight matrix for $f_3(\cdot)$ could provide faster convergence speed and better system performance. The weighted deep stream feature is computed as follows:

$$\tilde{X}_w = W_{\text{inter}}f_3(\tilde{X}_s) + f_2(\tilde{X}_d) \quad (2)$$

where $f_3(\cdot)$ has the same definition of $f_1(\cdot)$. The final combined feature is defined as:

$$\tilde{X}_{s,d} = \text{concat}(\tilde{X}_w, \tilde{X}_s) \quad (3)$$

The proposed feature fusion strategy shown in Figure 2 can transmit the speaker-related information into the shallow features and thus enhance the feature diversity.

2.4. Multi-head self-attention

The core of Transformer [9] is utilizing multi-head self-attention mechanism to learn the long-time sequence information. Multi-head self-attention is designed to integrate the information from different representation subspaces at different positions [9]. Each attention head uses the scaled dot-product attention to map a query and a set of key-value pairs to an output. Assumed $Q \in \mathbb{R}^{t_q \times d_q}$ denotes the query matrix, $K \in \mathbb{R}^{t_k \times d_k}$ denotes the key matrix and $V \in \mathbb{R}^{t_v \times d_v}$ denotes the value projection matrix, multi-head self-attention can be computed as:

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O \quad (4)$$

$$\text{head}_i = \text{SelfAttn}(QW_i^Q, KW_i^K, VW_i^V) \quad (5)$$

$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (6)$$

where the scalar $1/\sqrt{d_k}$ is the normalization term, head_i is the i^{th} head of multi-head self-attention, W^O , W_i^Q , W_i^K and W_i^V are learnable parameters.

2.5. Joint CTC/Attention Loss

The joint CTC/Attention architecture [20] profits from both CTC and attention-based models. The objective function is a logarithmic linear combination of the CTC and attention-based objective functions,

$$\mathcal{L} = \lambda \log P_{\text{ctc}}(C|\tilde{X}_{s,d}) + (1 - \lambda)P_{\text{att}}(C|\tilde{X}_{s,d}) \quad (7)$$

where P_{ctc} represents CTC loss [21], P_{att} represents attention loss, λ is a trade-off hyperparameter among CTC and attention loss, which satisfying $0 \leq \lambda \leq 1$.

3. Experiments

3.1. Experimental setups

Experiments were conducted on HKUST Mandarin telephone [17] and Librispeech [18] speech recognition tasks. The

HKUST mandarin telephone speech corpus contains about 170 hours training data while 5 hours for the development set and 5 hours for the evaluation set. The size of the training set of Librispeech corpus is 1000 hours and we used two standard test sets [18] (‘‘test-clean’’ and ‘‘test-other’’) for the evaluation.

For both the high-resolution and low-resolution spectrograms, we used 129-dimensional features, including 120-dimensional filterbanks (static, Δ and $\Delta\Delta$) and 9-dimensional Kaldi pitch (static, Δ and $\Delta\Delta$). The features were computed within a 25 ms window and the window shift was 10 ms. The size of the high-resolution spectrogram input X_d was 4 times larger than the size of the low-resolution spectrogram input X_s by rescaling operation. In addition, we followed the ESPNet [22] recipe to preprocess both corpora.

Our experiments were based on the existing recipes in the end-to-end speech processing toolkit ESPNet [22]. During the network training, $d_m = 256$ in equation (1) was used inside Transformer ASR model. For the joint CTC/Attention loss function, $\lambda = 0.3$ was used. We trained 20 epochs for the HKUST mandarin telephone speech recognition task and 120 epochs for the Librispeech speech recognition task.

3.2. Experimental results

Table 2: Character error rate (CER) of the Transformer ASR systems only using the shallow or deep stream features, or using both features on the HKUST test set.

Shallow Stream	Deep Stream	CER
✓		23.5
	✓	58.7
✓	✓	21.7

In order to explore the impact of the shallow stream features or deep stream features on the final system performance, we conducted experiments on the Transformer ASR systems with the multi-level acoustic feature extraction module shown in Figure 1. Table 2 shows the performance of the Transformer ASR systems only using the shallow or deep streams, or using both features on the HKUST test set. When only using the deep stream, the character error rate (CER) was up to 58.7% and it was too weak to use as the output receptive field of the deep stream. This result might be caused by losing lots of acoustic detailed information for distinguishing different phones or words. Moreover, when using the deep stream features as a complementary part for the shallow stream features, the resulted Transformer ASR system outperformed the baseline only using shallow stream features by a CER reduction of 8.3% relative on the HKUST test set.

To further explore the hidden information of the shallow and deep branches in the multi-level acoustic feature extraction module, we used t-SNE [23] projection to visualize the distribution of the utterance-level shallow and deep stream features. As shown in Figure 3, since the distance between two arbitrary feature points among the same category is small while the center point of two arbitrary different categories have a far way feature distance, the deep stream features could be used to identify different speakers. It indicates that such features contain the inherent characteristics of vocal cords of different speakers, which are missing in the shallow stream features and might be helpful for speech recognition tasks.

Table 3 shows the CER performance of the Transformer ASR systems using different fusion strategies on the HKUST

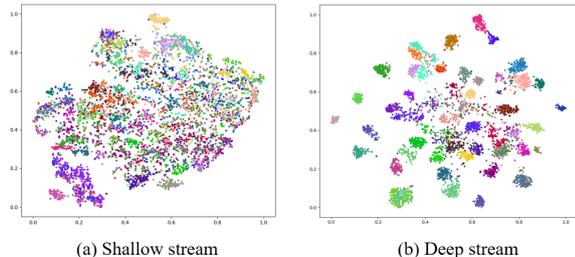


Figure 3: The t-SNE [23] visualization for the (a) shallow stream features and (b) deep stream features on the HKUST test set: different colors represent different speakers. T-SNE with default settings in the sklearn package [24] is applied for both streams.

Table 3: CER of the Transformer based ASR systems using different fusion strategies on the HKUST test set.

Fusion strategy	CER
Concatenation	22.2
Addition	21.9
FCF(Ours)	21.7

test set. As shown in Table 3, compared with the ASR systems using the “concatenation” and “fusion” strategies, the Transformer ASR system using the proposed FCF fusion strategy had the best performance. This indicated that the FCF fusion strategy across the frequency and time axes might be more suitable for capturing the hidden complementary information between the shallow and deep stream features.

Table 4: The effect of the number of the bottleneck residual blocks used on the deep branch in the multi-level acoustic feature extraction module on the HKUST Mandarin telephone speech recognition task.

#blocks	development set	test set
4	22.4	22.1
5	22.3	21.9
6	22.2	21.7

Table 4 shows the Transformer ASR system performance with respect to the number of the bottleneck residual blocks used on the deep branch on the HKUST Mandarin telephone speech recognition task. The number of the bottleneck residual blocks ranged from 4 to 6. As the number of the bottleneck residual blocks increases, the corresponding CER decreases on both the HKUST development and test sets. It would be suggested that more robust speaker-related information can be obtained at a deeper layer.

We also compared the proposed Transformer ASR system with other state-of-the-art ASR systems on the HKUST Mandarin telephone speech recognition task. As shown in Table 5, compared with other ASR systems, CER reductions of 1.4%-2.4% absolute (6%-10% relative) were obtained from the proposed speech recognition system on the HKUST corpus.

We also evaluated our framework on the Librispeech speech recognition task. As shown in Table 6, the proposed Transformer ASR system using the multi-level acoustic feature extraction module can achieve state-of-the-art performance on the Librispeech speech recognition task. Especially, compared with

Table 5: Comparison with different state-of-the-art ASR systems on the HKUST Mandarin telephone speech recognition task.

Model	CER
Chain-TDNN [25]	23.7
Self-attention Aligner [26]	24.1
CIF [4]	23.1
D-Att [27]	23.3
Ours	21.7

Table 6: Experiments in word error rate (WER) on the Librispeech speech recognition task. The ESPNet [22] was reimplemented with the setting $adim=256$, which is comparable with the proposed Transformer based ASR system with $d_m = 256$.

ASR systems	test-clean	test-other
RWTH [8]	3.8	8.8
DEJA-VU [28]	2.9	6.7
MEL-t-Fusion-Late [29]	3.3	7.2
ESPNet [22]	3.1	6.4
Ours	2.5	5.8

ESPNet, our final ASR system gave the WER reductions of 19.4% relative and 9.4% relative on the “test-clean” and “test-other” sets respectively.

4. Conclusions

This paper proposed a novel multi-level acoustic feature extraction framework that can be easily applied into Transformer based models. The proposed framework contained a shallow stream with less network layers and a deep stream with more network layers. Utterance-level speaker-invariant information derived from the deep stream with low-resolution spectrograms was used to improve the feature diversity, and then a novel feature correlation based fusion strategy was used to integrate the shallow and deep stream features across the frequency and time domains. Experimental results showed that the Transformer based ASR system using the fused features gave state-of-the-art performance on both HKUST Mandarin telephone and Librispeech speech recognition tasks. Future works will focus on applying the proposed feature extraction framework into other ASR models.

5. Acknowledgements

This work was supported in part by the National Key R&D Program of China (2020YFC2004100), National Natural Science Foundation of China (NSFC U1736202, NSFC 61771461, NSFC 62106255), Shenzhen Peacock Team Project (Grant No. KQTD20200820113106007) and Shenzhen Science and Technology Program (Grant No. JCYJ20210324115810030).

6. References

- [1] D. Amodei, S. Ananthanarayanan, R. Anubhai, J. Bai, E. Battenberg, C. Case, J. Casper, B. Catanzaro, Q. Cheng, G. Chen *et al.*, “Deep speech 2: End-to-end speech recognition in english and mandarin,” in *Proceedings of International Conference on Machine Learning (ICML)*, 2016, pp. 173–182.
- [2] X. Chen, Y. Wu, Z. Wang, S. Liu, and J. Li, “Developing real-time streaming transformer transducer for speech recognition on large-

- scale dataset,” in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021, pp. 5904–5908.
- [3] C.-C. Chiu, T. N. Sainath, Y. Wu, R. Prabhavalkar, P. Nguyen, Z. Chen, A. Kannan, R. J. Weiss, K. Rao, E. Gonina *et al.*, “State-of-the-art speech recognition with sequence-to-sequence models,” in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2018, pp. 4774–4778.
 - [4] L. Dong and B. Xu, “Cif: Continuous integrate-and-fire for end-to-end speech recognition,” in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 6079–6083.
 - [5] Z. Meng, S. Parthasarathy, E. Sun, Y. Gaur, N. Kanda, L. Lu, X. Chen, R. Zhao, J. Li, and Y. Gong, “Internal language model estimation for domain-adaptive end-to-end speech recognition,” in *Proceedings of IEEE Spoken Language Technology Workshop (SLT)*, 2021, pp. 243–250.
 - [6] H. Miao, G. Cheng, C. Gao, P. Zhang, and Y. Yan, “Transformer-based online ctc/attention end-to-end speech recognition architecture,” in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 6084–6088.
 - [7] J. Chorowski, D. Bahdanau, K. Cho, and Y. Bengio, “End-to-end continuous speech recognition using attention-based recurrent NN: First results,” *arXiv preprint arXiv:1412.1602*, 2014.
 - [8] C. Lüscher, E. Beck, K. Irie, M. Kitzka, W. Michel, A. Zeyer, R. Schlüter, and H. Ney, “Rwth asr systems for librispeech: Hybrid vs attention–w/o data augmentation,” *arXiv preprint arXiv:1905.03072*, 2019.
 - [9] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Proceedings of Conference on Neural Information Processing Systems (NIPS)*, 2017, pp. 5998–6008.
 - [10] L. Dong, S. Xu, and B. Xu, “Speech-transformer: a no-recurrence sequence-to-sequence model for speech recognition,” in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2018, pp. 5884–5888.
 - [11] H. Luo, S. Zhang, M. Lei, and L. Xie, “Simplified self-attention for transformer-based end-to-end speech recognition,” in *Proceedings of IEEE Spoken Language Technology Workshop (SLT)*, 2021, pp. 75–81.
 - [12] T. Nakatani, “Improving transformer-based end-to-end speech recognition with connectionist temporal classification and language model integration,” in *Proceedings of ISCA INTERSPEECH*, 2019, pp. 1480–1412.
 - [13] X. Wang, R. Li, S. H. Mallidi, T. Hori, S. Watanabe, and H. Hermansky, “Stream attention-based multi-array end-to-end speech recognition,” in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 7105–7109.
 - [14] F.-J. Chang, M. Radfar, A. Mouchtaris, B. King, and S. Kunzmann, “End-to-end multi-channel transformer for speech recognition,” in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021, pp. 5884–5888.
 - [15] D. T. Toledano, M. P. Fernández-Gallego, and A. Lozano-Diez, “Multi-resolution speech analysis for automatic speech recognition using deep neural networks: Experiments on TIMIT,” *PLoS one*, vol. 13, no. 10, pp. 1–24, 2018.
 - [16] X. Wang, R. Girshick, A. Gupta, and K. He, “Non-local neural networks,” in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 7794–7803.
 - [17] Y. Liu, P. Fung, Y. Yang, C. Cieri, S. Huang, and D. Graff, “Hkust/mts: A very large scale mandarin telephone speech corpus,” in *Proceedings of IEEE International Symposium on Chinese Spoken Language Processing (ISCSLP)*, 2006, pp. 724–735.
 - [18] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, “Librispeech: an asr corpus based on public domain audio books,” in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2015, pp. 5206–5210.
 - [19] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, “Mobilenetv2: Inverted residuals and linear bottlenecks,” in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 4510–4520.
 - [20] S. Kim, T. Hori, and S. Watanabe, “Joint ctc-attention based end-to-end speech recognition using multi-task learning,” in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2017, pp. 4835–4839.
 - [21] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, “Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks,” in *Proceedings of International Conference on Machine Learning (ICML)*, 2006, pp. 369–376.
 - [22] S. Watanabe, T. Hori, S. Karita, T. Hayashi, J. Nishitoba, Y. Unno, N. E. Y. Soplin, J. Heymann, M. Wiesner, N. Chen *et al.*, “Espnet: End-to-end speech processing toolkit,” *arXiv preprint arXiv:1804.00015*, 2018.
 - [23] L. Van Der Maaten, “Accelerating t-sne using tree-based algorithms,” *The journal of machine learning research*, vol. 15, no. 1, pp. 3221–3245, 2014.
 - [24] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg *et al.*, “Scikit-learn: Machine learning in python,” *The Journal of machine Learning research*, vol. 12, pp. 2825–2830, 2011.
 - [25] D. Povey, V. Peddinti, D. Galvez, P. Ghahremani, V. Manohar, X. Na, Y. Wang, and S. Khudanpur, “Purely sequence-trained neural networks for ASR based on lattice-free MMI,” in *Proceedings of ISCA INTERSPEECH*, 2016, pp. 2751–2755.
 - [26] L. Dong, F. Wang, and B. Xu, “Self-attention aligner: A latency-control end-to-end model for ASR using self-attention network and chunk-hopping,” in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 5656–5660.
 - [27] C. Gao, G. Cheng, J. Zhou, P. Zhang, and Y. Yan, “Non-autoregressive deliberation-attention based end-to-end ASR,” in *Proceedings of IEEE International Symposium on Chinese Spoken Language Processing (ISCSLP)*, 2021, pp. 1–5.
 - [28] A. Tjandra, C. Liu, F. Zhang, X. Zhang, Y. Wang, G. Synnaeve, S. Nakamura, and G. Zweig, “Deja-vu: Double feature presentation and iterated loss in deep transformer networks,” in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 6899–6903.
 - [29] T. Lohrenz, Z. Li, and T. Fingscheidt, “Multi-encoder learning and stream fusion for transformer-based end-to-end automatic speech recognition,” *arXiv preprint arXiv:2104.00120*, 2021.