
Mini-Project

An Informal Review of Relevant
Literature

Robert Flynn

wav2vec 2.0: A Framework for Self-Supervised Learning of Speech Representations (Baevski et al., 2020)

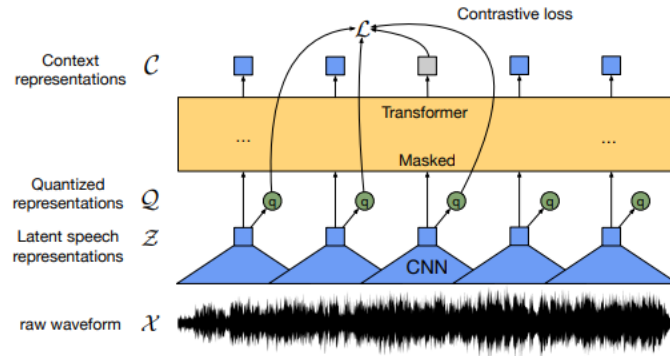


Figure 1: Diagram detailing the contrastive self-supervised pre-training for Wav2vec 2.0

- This model uses a self-supervised learning framework to build acoustic representations from unlabeled speech
- The model's architecture consists of a feature encoder that uses a convolutional architecture, and a transformer architecture which intakes the features and outputs contextual speech representations
- For self-supervised pre-training the model employs a contrastive learning approach. A subset of the latent speech features representations are masked, and the model is trained with the objective of identifying the correct latent quantized representation when marginalizing over a set of “distractors” sampled from the sequence
- Figure 1 depicts the pre-training task. This approach is inspired by the masked language model pre-training used by the BERT model (Devlin et al., 2018)
- Unlabeled speech data is often readily available. This work shows that once good acoustic representations have been learnt, the model can be fine-tuned to transcribe text with a small amount of labeled data. More labeled data is always better, however this approach shows reasonable results can be obtained when fine-tuning on as little of 10 minutes of labelled audio: “when using only 10 minutes of labeled data, our approach achieves word error rate (WER) 4.8/8.2 on the clean/other test sets of Librispeech”
- For fine-tuning, a linear-projection layer is initialized on top of the (transformer) context layer. This layer is trained to predict a set of classes, which consists of the vocabulary of the target text, for librispeech they use 29 tokens, to denote characters and word boundaries.
- For fine-tuning the model uses Connectionist Temporal Classification (CTC) (Graves et al., 2006) as its loss function. An intuitive explanation of this process is provided by Hannun (2017).

- This is an encoder only model with the linear projection layer predicting characters based on the acoustic representations at each time-step; consequently, the probabilities of the outputted characters are conditionally independent. Language models are often used with such architectures to find a more probable character sequence given the outputted logits.

ROBUST WAV2VEC 2.0: ANALYZING DOMAIN SHIFT IN SELF-SUPERVISED PRE-TRAINING (Hsu et al., 2021)

- The Wav2vec2 paper showed that once good acoustic representations have been learnt, models can be fine-tuned to transcribe audio effectively with a much smaller amount of labelled data than if no pre-training was performed.
- Pre-training for Wav2vec2 took place on audio in the same domain as the audio for fine-tuning, and evaluation. Often there will be situations where there is not enough transcribed audio from the desired target domain to fine-tune a model such as Wav2vec2.
- This paper examines whether pre-training on in-domain data then fine-tuning on Out-Of-Domain (OOD) data is beneficial for producing an ASR system that is effective in the target domain

“In the supervised learning paradigm, practitioners who would like to build a system for a new domain, can either train on existing OOD labeled data or build a corpus of labeled data in the new domain. With pre-training, we have a third option: collect unlabeled data in the new domain and fine-tune on existing labeled OOD data. This has the clear advantage of unlabeled in-domain being often much easier to obtain than transcribed in-domain data.”

X	<i>TED-LIUM (TD) dev WER</i>					
	FT on TD-10h		FT on LS-10h		FT on SB-10h	
	PT on X	X+TD	PT on X	X+TD	PT on X	X+TD
None	diverge	9.93	diverge	10.99	diverge	11.32
SF	12.12	9.60	14.82	11.08	99.63	11.04
LS	9.81	8.59	12.92	8.91	13.08	10.39
SF+LS	9.13	8.91	10.61	9.67	12.25	10.75

X	<i>LibriSpeech (LS) dev-other WER</i>					
	FT on TD-10h		FT on LS-10h		FT on SB-10h	
	PT on X	X+LS	PT on X	X+LS	PT on X	X+LS
None	diverge	14.60	diverge	10.53	diverge	17.92
SF	28.91	14.30	20.36	10.44	94.38	15.53
TD	23.44	12.81	15.36	9.71	27.50	15.46
SF+TD	20.50	13.58	14.42	10.39	21.99	13.89

X	<i>Switchboard (SB) RT03 WER</i>					
	FT on TD-10h		FT on LS-10h		FT on SB-10h	
	PT on X	X+SF	PT on X	X+SF	PT on X	X+SF
None	diverge	18.90	diverge	19.30	diverge	10.80
TD	35.70	16.20	34.60	17.40	18.70	11.00
LS	33.60	17.80	36.50	16.10	18.20	11.00
TD+LS	29.70	17.40	28.90	16.90	15.60	10.80

Figure 2: Validation WER on TD, LS, and SB of models pre-trained (PT) on various subsets of TD, LS, SB, and fine-tuned (FT) on TD-10h, LS-10h, or SB-10h.

- The table in figure 2 presents the key findings of this work. The main takeaway from these results is that pre-training on data for the target domain can lead to sizeable improvements in WER, despite only fine-tuning on OOD data.
- Further results from this work show that pre-training on multiple domains improves the models ability to perform on previously unseen domains at evaluation time.

What are the implications of this work with regard to the mini-project?

It is likely that we will be able to obtain a relatively large amount of unlabelled data for our target domain of oral history. While it may be possible to acquire some data that shares similar characteristics to our target domain, it may not be an optimal amount to train/fine-tune a large ASR model on its own. Work from this paper suggests that pre-training an acoustic model with data that includes large amounts of audio from our target domain, then fine-tuning on a combination of: any near-target domain transcribed audio we can acquire, and on other datasets commonly used in ASR (librespeech e.t.c.), may yield good results.

Our current baseline model uses a Robust Wav2vec2 model that has been finetuned on librespeech¹. A version of Robust Wav2vec2 that has not been finetuned is also available on the hugging face library², this model has been pre-trained in an unsupervised fashion on a variety of commonly used corpora (libri light, Common voice, Switchboard and Fisher). We could potentially use this model as the starting point for further pre-training on Oral History unlabeled data.

Effective Sentence Scoring Method Using BERT for Speech Recognition (Shin et al., 2019)

- **A given sentence:**
`move the vat over the hot fire`
- **A set of instances we create:**
 1. Input = [MASK] the vat over the hot fire
Label = move
 2. Input = move [MASK] vat over the hot fire
Label = the
 - ...
 7. Input = move the vat over the hot [MASK]
Label = fire

Figure 3: A depiction of the BERT model being used to predict word probabilities, when each word in the sentence is masked iteratively

¹<https://huggingface.co/facebook/wav2vec2-large-robust-ft-libri-960h>

²<https://huggingface.co/facebook/wav2vec2-large-robust>

- Language Models (LMs) are commonly incorporated into ASR acoustic models to improve performance. These models can either be implicitly integrated in a fully sequence-to-sequence (seq-to-seq) fashion, or utilized through beam search decoding in CTC encoder only models.
- This paper uses the Listen, Attend and Spell (Chan et al., 2015) ASR model, which employs an attention based BiLSTM type architecture. In this paper the LAS model was trained with a CTC objective, as slightly more recent work demonstrated that the left-to-right constraints of CTC help the model learn speech text alignments (Hori et al., 2017). This is because regular attention mechanisms do not take advantage of the monotonic alignment between the inputted audio and outputted text (the order of words in the outputted text should always be given in the same order in which they are spoken) (Lugosch, 2020).
- LMs used for decoding and beam search re-scoring in ASR models are often unidirectional, processing the input sequence from left-to-right. Bi-directional models offer the advantage of being able to assess the probability of a word in the context of the entire sentence, not just the words prior.
- Bi(directional)-LSTM LMs have been previously applied to ASR however “there is no interaction between the past and the future words in the biLMs” as the forward and backward representations are not fused
- The BERT (Devlin et al., 2018) architecture combats these issues. BERT model is also trained on extremely large corpus, with multi-headed attention blocks, allowing the model to better capture long-range token level dependencies (it can accurately model the relationships between words that are far away from each other)
- This paper looks at using the BERT model to re-rank candidate transcribed sentences. Their technique utilizes BERT’s Pre-training task, iteratively masking each word in a sentence, and summing the probabilities of each target to attain an overall “likelihood” of the sentences. This process is shown in figure 3
- This score is linearly combined with the prior score attributed to the sentences by the LAS model, with some weighting attributed to each probability.
- These new scores are then used to re-rank the candidate sentence and attain a new (hopefully) more likely sentence. In this paper they trained the BERT model on the libre-speech corpus, along with a range of other uni-directional LMs, with their proposed methodology showing the best performance.

What are the implications of this work with regard to the mini-project?

A similar technique to this could be applied to any ASR system to find more likely word sequences. Although they train BERT on the Librespeech corpus this may not be necessary, and simply using a pre-trained BERT model would likely be sufficient to see improvements. Utilizing large attention based LMs such as BERT should be advantageous for our project, where deciphering the speaker’s utterances often requires disambiguation through examining the context.

Efficiently Fusing Pretrained Acoustic and Linguistic Encoders for Low-resource Speech Recognition (Yi et al., 2021)

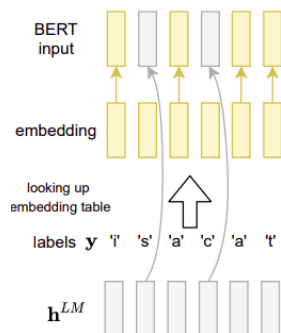


Figure 4: Depiction of acoustic representations being integrated into a BERT encoder. Acoustic representations (h^{AC}) are passed through a fully connected layer, mapping $h^{AC} \rightarrow h^{LM}$. These representations (h^{LM}) are randomly mixed into the BERT encoder during training

- This work looks at fusing acoustic encoders (Wav2vec2) with linguistic encoders (Bert), for speech domains where there exists little labeled data
- The authors categorize ASR approaches into two sets: Pipeline methods, that use separated acoustic and linguistic models; and, end-to-end models that integrate all components into one, this generally involves a seq-to-seq model with an encoder-decoder framework that intakes speech and outputs text.
- These encoder-decoder frameworks have the capacity to outperform pipeline methods on most public datasets. However, these networks require large amounts (100s of hours) of transcribed speech to perform effectively on the target domain, which is problematic for low-resource ASR.
 - Transfer learning (where knowledge learned from other tasks is applied to the target one) is one method of adapting these models to low-resource domains, this *generally* (See SpeechStew (Chan et al., 2021)) requires “domain-similar” labelled data to pre-train the model on, which may be hard to find.
 - Another method (Jiang et al., 2019) is to train the acoustic encoder in an unsupervised fashion (masked predictive coding) on unlabeled data, then add the decoder to the model while fine-tuning on labeled data. However, the decoder cannot be pre-trained separately as it relies on the acoustic representations from the encoder. Consequently, the decoder is only able to learn a model of language during labelled fine-tuning.
- Pipeline methods (including wav2vec2) require much less labelled data due to the effectiveness of self-supervised pre-training techniques, however pipeline components (i.e the language model) are often combined through some fixed weighting, which is “inflexible”
- Because of these factors the authors choose to explore better methods of combining linguistic encoders (BERT) with acoustic models. “The fused model has been separately exposed

to adequate speech and text data, so that it only needs to learn the transfer from speech to language during fine-tuning with limited labeled data.”

- During fine-tuning a random selection of wav2vec2’s Acoustic representations are randomly mapped into BERT’s input layers. These two networks are connected through a CIF (continuous-integrate-and-fire) (Dong and Xu, 2019) mechanism, which is a variant of Attention, designed to combat its monotonic alignment issues.
- Figure 4 depicts how the encoders are fused during fine-tuning. Initially BERT will not know how to interpret these representations (h^{LM}), and will “view” them as masked input, mainly relying on the transcription labels to make predictions. Through fine-tuning BERT is able to “learn” the meaning of the h^{LM} representations, and utilize them along with prior linguistic knowledge to make predictions.

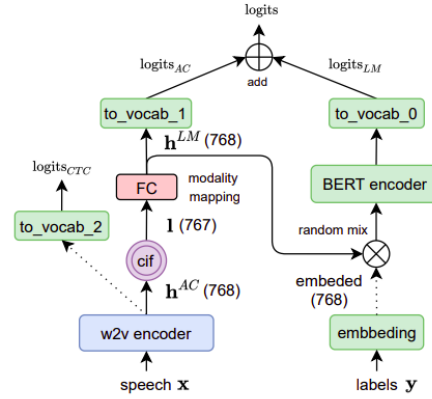


Figure 5: Depictions of the entire fusion network. Dotted lines represent modules that are ignored during inference

- Figure 5 presents the entire network. Logits from both encoders are combined and the entire network is trained through cross-entropy loss. The Wav2vec2 encoder receives additional supervision from its regular CTC task.

What are the implications of this work with regard to the mini-project?

This approach shows reasonable improvements on some low-resource data-sets compared to a selection of previous work. The authors provide good justification for their approach of acoustic and LM encoder fusion for low-resource settings. More recent work (Zheng et al., 2021) provides a direct comparison with the previously mentioned BERT re-scoring (Shin et al., 2019) system on a (different) range of low-resource languages. Results showed that Yi et al. (2021)’s fusion approach performed worse than BERT re-scoring method in all settings. Catastrophic forgetting (where knowledge learnt during pre-training is lost during fine-tuning) is likely the reason for this models comparatively poor performance. Fusing the Wav2vec2 directly with BERT’s input may cause linguistic knowledge learnt via pre-training to be forgotten (Zheng et al., 2021).

Selecting a simpler and more effective approach such as the BERT re-scoring method would be best for our project. The Wav-BERT (Zheng et al., 2021) system fuses both models (wav2vec2

+ BERT) through a “representation aggregation module” and a Gated Attention operation (Xue et al., 2019; Zhang et al., 2018), that helps overcome the catastrophic forgetting issue. Wav-BERT shows improved results over previous models on a selection of low-resource languages. The improvements seen from Wav-BERT are still only marginal when compared against the much simpler re-scoring (Shin et al., 2019) methods; additionally, the implementation Wav-BERT is likely outside the scope of our project.

Conformer: Convolution-augmented Transformer for Speech Recognition (Gulati et al., 2020)

- Transformers have enjoyed success on a range of sequence modelling tasks due their ability capture long range dependencies. However these networks are less capable of extracting “fine-grained local feature patterns”.
- Convolutional networks excel at capturing local information, hence their uptake in computer vision, however they require many layers and parameters to capture a sup-optimal global context.
- The authors: “hypothesize that both global and local interactions are important for being parameter efficient. To achieve this, we propose a novel combination of self-attention and convolution will achieve the best of both worlds”
- The paper introduces the Conformer network that utilizes convolutions and attention to produce their encoder, and an LSTM based decoder
- The Conformer improves on previous work including LAS (Chan et al., 2015) and Transformer-Transducers (Zhang et al., 2020a) networks on the Librespeech dataset
- This architecture utilizes SpecAugment (Park et al., 2019), which is method used for data augmentation in ASR which has led to improved performance on models such as LAS (Chan et al., 2015).

What are the implications of this work with regard to the mini-project?

This version of the Conformer network has been super-seeded by more recent models including Wav2vec2. Additionally, the approach used in this paper only utilizes labeled data, which may make it harder to adapt to low-resource settings. The author’s investigate different parameter sizes with this model and see reasonable results with model sizes as low as 10M parameters, training a low parameter model should require less data and resources.

This methodology presents good results, and it is worth being aware of this architecture and paper as later literature builds upon this. A Wav2vec2 based approach should still be more suited to our project, however approaches involving this architecture may be worth considering. Wav2vec2 offers the clear advantage of unsupervised pre-training but requires additional methods to incorporate linguistic models for optimal performance.

SpecAugment: A Simple Data Augmentation Method for Automatic Speech Recognition (Park et al., 2019)

- Data Augmentation (DA) has been applied in state-of-the-art systems for computer vision tasks, and has also seen success in the speech domain. DA can improve the robustness of

models, and can be easily implemented for domains that use continuous data.

- SpecAugment can be directly applied to the feature inputs of a model. The authors select deformations for DA that help ASR networks learn useful features, and improve robustness to: temporal deformations, partial loss of frequency information, and partial loss of speech segments. The following mentions each of the deformations performed, see the paper for more details:
 - Time Warping
 - Frequency Masking
 - Time Masking
- Results from this work show that the SpecAugment DA method can improve the WER of a network such as LAS Chan et al. (2015). Their approach shows sizeable improvements to the original networks performance, outperforming **prior** work.
- The authors note that Time Warping, although helpful, is not a major factor in the performance improvements, and this deformation can be dropped given any budgetary/hardware limitations.
- This DA method can cause networks that previously overfitted on a given dataset, to underfit. This means deeper networks can be trained on smaller datasets using this approach.

What are the implications of this work with regard to the mini-project?

DA, and particularly the SpecAugment method, may be helpful for our project. Such techniques may increase the robustness of our model, and increase/augment the amount of data we have. SpecAugment is integrated into the current SOTA on libraspeech (Zhang et al., 2020b). These techniques would be trivial method to implement as there exists an implementation of this method on GitHub³ which allows for integration with Pytorch and Tensorflow based models.

³<https://github.com/DemisEom/SpecAugment>

Pseudo-Labeling

This section will overview a selection of papers that explore the use of pseudo-labelling for the task of semi-supervised ASR.

SELF-TRAINING FOR END-TO-END SPEECH RECOGNITION (Kahn et al., 2019)

Pseudo labelling is a technique that utilizes a trained model to provide labels for unlabeled data, this process works as follows:

1. We have a labelled dataset D , unlabelled audio X , and unpaired text data Y
2. Train an acoustic model on paired labelled D
3. Train a language model on unpaired text data Y
4. Combine the acoustic model and the language model and generate psuedo-labels for unlabelled audio X , creating new dataset \bar{D}
5. Train new acoustic model on a combination of D and \bar{D}

The authors employ a heuristic based filtering method, and a log likelihood based confidence score, to remove undesirable examples from \bar{D} . Model ensembles can be used to decrease the impact of sample noise (inaccurate psuedo-labels). This can be performed using the process described in the above list, using N randomly initialized models, producing a collection: $\{\bar{D}_1, \dots, \bar{D}_N\}$. When training the final acoustic model psuedo-labels are uniformly sampled from one of the N models each epoch.

Future Reading/Papers To Be Added:

- A COMPARISON OF TECHNIQUES FOR LANGUAGE MODEL INTEGRATION IN ENCODER-DECODER SPEECH RECOGNITION (Toshniwal et al., 2018)
- Pushing the Limits of Semi-Supervised Learning for Automatic Speech Recognition (Zhang et al., 2020b)
- END-TO-END ASR: FROM SUPERVISED TO SEMI-SUPERVISED LEARNING WITH MODERN ARCHITECTURES (Synnaeve et al., 2019)
- SELF-TRAINING FOR END-TO-END SPEECH RECOGNITION (Kahn et al., 2019)
- Improved Noisy Student Training for Automatic Speech Recognition (Park et al., 2020)
- Iterative Pseudo-Labeling for Speech Recognition (Xu et al., 2020)

References

- Alexei Baevski, Henry Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. *CoRR* abs/2006.11477.
- William Chan, Navdeep Jaitly, Quoc V. Le, and Oriol Vinyals. 2015. Listen, attend and spell. *CoRR* abs/1508.01211.
- William Chan, Daniel S. Park, Chris Lee, Yu Zhang, Quoc V. Le, and Mohammad Norouzi. 2021. Speechstew: Simply mix all available speech recognition data to train one large neural network. *CoRR* abs/2104.02133.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR* abs/1810.04805.
- Linhao Dong and Bo Xu. 2019. CIF: continuous integrate-and-fire for end-to-end speech recognition. *CoRR* abs/1905.11235.
- Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber. 2006. Connectionist temporal classification: Labelling unsegmented sequence data with recurrent neural networks. In *Proceedings of the 23rd International Conference on Machine Learning*. Association for Computing Machinery, New York, NY, USA, ICML '06, page 369–376.
- Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, and Ruoming Pang. 2020. Conformer: Convolution-augmented transformer for speech recognition.
- Awni Hannun. 2017. Sequence modeling with ctc. *Distill* <https://distill.pub/2017/ctc>.
- Takaaki Hori, Shinji Watanabe, Yu Zhang, and William Chan. 2017. Advances in joint ctc-attention based end-to-end speech recognition with a deep CNN encoder and RNN-LM. *CoRR* abs/1706.02737.
- Wei-Ning Hsu, Anuroop Sriram, Alexei Baevski, Tatiana Likhomanenko, Qiantong Xu, Vineel Pratap, Jacob Kahn, Ann Lee, Ronan Collobert, Gabriel Synnaeve, and Michael Auli. 2021. Robust wav2vec 2.0: Analyzing domain shift in self-supervised pre-training. *CoRR* abs/2104.01027.
- Dongwei Jiang, Xiaoning Lei, Wubo Li, Ne Luo, Yuxuan Hu, Wei Zou, and Xiangang Li. 2019. Improving transformer-based speech recognition using unsupervised pre-training. *CoRR* abs/1910.09932.
- Jacob Kahn, Ann Lee, and Awni Hannun. 2019. Self-training for end-to-end speech recognition. *CoRR* abs/1909.09116.
- Loren Lugosch. 2020. Sequence-to-sequence learning with transducers.
- Daniel S. Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D. Cubuk, and Quoc V. Le. 2019. Specaugment: A simple data augmentation method for automatic speech recognition. *Interspeech 2019*.
- Daniel S. Park, Yu Zhang, Ye Jia, Wei Han, Chung-Cheng Chiu, Bo Li, Yonghui Wu, and Quoc V. Le. 2020. Improved noisy student training for automatic speech recognition. *Interspeech 2020*.

- Joonbo Shin, Yoonhyung Lee, and Kyomin Jung. 2019. Effective sentence scoring method using bert for speech recognition. In Wee Sun Lee and Taiji Suzuki, editors, *Proceedings of The Eleventh Asian Conference on Machine Learning*. PMLR, volume 101 of *Proceedings of Machine Learning Research*, pages 1081–1093.
- Gabriel Synnaeve, Qiantong Xu, Jacob Kahn, Edouard Grave, Tatiana Likhomanenko, Vineel Pratap, Anuroop Sriram, Vitaliy Liptchinsky, and Ronan Collobert. 2019. End-to-end ASR: from supervised to semi-supervised learning with modern architectures. *CoRR* abs/1911.08460.
- Shubham Toshniwal, Anjuli Kannan, Chung-Cheng Chiu, Yonghui Wu, Tara N Sainath, and Karen Livescu. 2018. A comparison of techniques for language model integration in encoder-decoder speech recognition.
- Qiantong Xu, Tatiana Likhomanenko, Jacob Kahn, Awni Hannun, Gabriel Synnaeve, and Ronan Collobert. 2020. Iterative pseudo-labeling for speech recognition.
- Lanqing Xue, Xiaopeng Li, and Nevin L. Zhang. 2019. Not all attention is needed: Gated attention network for sequence data. *CoRR* abs/1912.00349.
- Cheng Yi, Shiyu Zhou, and Bo Xu. 2021. Efficiently fusing pretrained acoustic and linguistic encoders for low-resource speech recognition. *CoRR* abs/2101.06699.
- Jiani Zhang, Xingjian Shi, Junyuan Xie, Hao Ma, Irwin King, and Dit-Yan Yeung. 2018. Gaan: Gated attention networks for learning on large and spatiotemporal graphs. *CoRR* abs/1803.07294.
- Qian Zhang, Han Lu, Hasim Sak, Anshuman Tripathi, Erik McDermott, Stephen Koo, and Shankar Kumar. 2020a. Transformer transducer: A streamable speech recognition model with transformer encoders and rnn-t loss.
- Yu Zhang, James Qin, Daniel S. Park, Wei Han, Chung-Cheng Chiu, Ruoming Pang, Quoc V. Le, and Yonghui Wu. 2020b. Pushing the limits of semi-supervised learning for automatic speech recognition.
- Guolin Zheng, Yubei Xiao, Ke Gong, Pan Zhou, Xiaodan Liang, and Liang Lin. 2021. Wav-bert: Cooperative acoustic and linguistic representation learning for low-resource speech recognition.