

UNSUPERVISED PRE-TRAINING OF BIDIRECTIONAL SPEECH ENCODERS VIA MASKED RECONSTRUCTION

Weiran Wang^{*}

Qingming Tang[†]

Karen Livescu[†]

^{*}Salesforce Research

weiran.wang@salesforce.com

[†]Toyota Technological Institute at Chicago

{qmtang, klivescu}@ttic.edu

ABSTRACT

We propose an approach for pre-training speech representations via a masked reconstruction loss. Our pre-trained encoder networks are bidirectional and can therefore be used directly in typical bidirectional speech recognition models. The pre-trained networks can then be fine-tuned on a smaller amount of supervised data for speech recognition. Experiments with this approach on the LibriSpeech and Wall Street Journal corpora show promising results. We find that the main factors that lead to speech recognition improvements are: masking segments of sufficient width in both time and frequency, pre-training on a much larger amount of unlabeled data than the labeled data, and domain adaptation when the unlabeled and labeled data come from different domains. The gain from pre-training is additive to that of supervised data augmentation.

Index Terms— Unsupervised representation learning, Pre-training, Masked reconstruction

1. INTRODUCTION

We study the problem of improving speech recognition via unsupervised pre-training, possibly on external data. Unsupervised pre-training has a long history in the field of speech recognition. Restricted Boltzmann Machines (RBMs) [1] were widely used to pre-train deep neural networks as part of a speech recognizer [2], often on the same transcribed data used for acoustic modeling. In recent years, however, RBM-based pre-training has been largely abandoned, because direct supervised training of deep neural networks has improved due to new techniques such as better initialization [3], non-saturating activation functions [4], and better control of generalization [5]. However, very recent work has begun to reconsider the value of unsupervised pre-training, specifically in the context of representation learning on a large set of unlabeled data, for use in supervised training on a smaller set of labeled data [6, 7, 8].

At the same time, in the area of natural language processing (NLP), unsupervised pre-trained representation learning has been extremely successful. In the past two years, several

approaches have been proposed for pre-trained text representations [9, 10, 11]. In particular, BERT [11] and its variants have enabled large improvements over the previous state of the art on a number of benchmark tasks [12].

In this paper we take inspiration from BERT-style pre-training, specifically its use of masked reconstruction loss, and adapt the idea for speech recognition. BERT is a bidirectional model that takes as input text that has had a certain percentage of randomly selected tokens masked, and attempts to reconstruct the masked text. The idea is that a model that can predict the missing data should provide a good representation of the important content. The same idea should hold for speech, but there are some significant differences between text and speech signals. In particular, the speech signal is continuous while text is discrete; and speech has much finer granularity than text, such that a single word typically spans a large sequence of contiguous frames. To handle these properties of speech, we take our second inspiration from recent work on speech data augmentation [13], which applies masks to the input in both the time and frequency domains. Thus, rather than randomly masking a certain percentage of frames (as in BERT training), we randomly mask some channels across all time steps of the input sequence, as well as contiguous segments in time. We experiment with a range of choices for the number and width of masks, and find that for appropriate choices our BERT-style pretraining significantly improves over strong speech recognition baselines.

2. RELATED WORK

Recent work has considered unsupervised learning for a variety of speech tasks. Some of this work is explicitly aimed at a “zero-speech” setting where no or almost no labeled data is available at all (e.g., [14, 15, 16, 17]), where the focus is to learn phonetic or word-like units, or representations that can distinguish among such units. Other work considers a variety of downstream supervised tasks, and some focuses explicitly on learning representations that generalize across tasks or across very different domains [6, 7, 18, 19]. This work uses a variety of training objectives, including autoencoder-based [15] and language model-like [7].

Specifically for our setting of unsupervised pre-training

Work done while Weiran Wang was at Amazon Alexa.

for supervised ASR, Schneider *et al.* [8] and Pascual *et al.* [6] learn unsupervised convolutional network-based representations, and show that they improve the performance of ASR trained on smaller labeled data sets. Their work relates to a number of other recent approaches for unsupervised representation learning [20, 21] based on the idea of maximizing (a lower bound on) mutual information (MI) between the current-time-step representation and future-time-step inputs (or shallow features of the inputs). Such approaches use either convolutional or unidirectional architectures to extract representations from audio, as their objective relies on the notion of “future”, which is not applicable for bidirectional models. These methods obtain impressive results, but are not directly applicable to pre-training bidirectional RNNs, though they can in principle be stacked with bidirectional RNNs. Concurrent work [22] combines a mutual information-based approach with vector quantization for learning discrete representations, which are then used as input to BERT—an example of stacking a bidirectional model on top of a unidirectional MI-based one.

Our work contrasts with prior work in several ways. First, to the best of our knowledge our work is the first to pre-train bidirectional RNNs for direct use in a speech recognizer and to show improved recognition in this setting. Besides the concurrent work of [22], we believe our work is also the first to use BERT-style masked reconstruction for representation learning for speech recognition. In addition, we use continuous spectrogram-based input, which allows us to explore both time- and frequency-domain masking, and produces an overall much simpler method. Finally, unlike other recent unsupervised pre-training approaches, we explicitly consider the problem of domain mismatch between the pre-training and fine-tuning data sets (see Section 4.3), and show that a simple adaptation layer can help address it.

3. PRE-TRAINING BY MASKED RECONSTRUCTION

The main idea of BERT training is to perturb the inputs by randomly masking tokens with some probability, and reconstruct the masked tokens at the output. Inspired by this idea, we perform representation learning for speech by masked reconstruction. Unlike the text domain where the inputs are discrete tokens, in the speech domain, the inputs are usually multi-dimensional feature vectors (e.g., energy in multiple frequency bands) in each frame, which are continuous and vary smoothly over time. Moreover, the time span of each frame is typically tens of milliseconds, much shorter than the span of the modeling unit in ASR. Our approach adapts the idea of masked reconstruction to the speech domain.

Our approach can also be viewed as extending the data augmentation technique SpecAugment [13], which was shown to be useful for supervised ASR, to unsupervised representation learning. We begin with a spectrogram representation of the input utterance. Viewing each input utterance

Fig. 1: Illustration of our masked reconstruction approach.

as an image of dimension $D \times T$, where D is the number of frequency bins and T the number of frames, we adopt the spectral masking technique of [13] for masking the inputs: We select m_F segments of the input in the frequency axis with random locations, whose widths are drawn uniformly from $\{0, 1, \dots, n_F\}$, and similarly select m_T segments in the time axis, with widths up to n_T , and set the selected pixels (time-frequency bins) to value 0. The intent is that masking in both frequency and time should encourage the network to exploit spatio-temporal patterns in the input.

Fig. 1 illustrates our approach. Each input utterance X is perturbed with a binary mask M of the same dimensions as X ($M_{td} = 0$ if X_{td} is being masked, for $t = 1, \dots, T$ and $d = 1, \dots, D$), and then passed through a feature extractor f consisting of several bidirectional recurrent neural layers followed by a linear layer, to obtain a high level representation (features) for each frame. Another (deep feedforward) network g is then used to reconstruct the input from the features. We measure the loss on the masked portion of the input:

$$\mathcal{L}(X, M; f, g) = \|(1 - M) \odot [X - g(f(M \odot X))]\|_{Fro}^2$$

where \odot denotes element-wise multiplication. Given a set of unlabeled utterances, we minimize the average of this reconstruction loss over the set. After unsupervised pre-training, we retain the LSTM layers of f and use them as initialization for supervised ASR training.

To augment the unlabeled data, we also make use of the speed perturbation method of [13], which performs linear interpolation along the time axis. Besides the original data, we use two additional speed factors 0.9 and 1.1, effectively obtaining 3 times as much data for pre-training.

4. EXPERIMENTS

4.1. Setup

We demonstrate our pre-training method using the LibriSpeech [23] and WSJ corpora¹. We explore a few settings with different amounts of data for unsupervised pre-training and supervised fine-tuning. Supervised training is always performed on WSJ, with either the *si84* partition (7040 utterances, 15 hours) or the *si284* partition (37.3K utterances, 80 hours) as the training set; the *dev93* partition (503 utterances) is used as development set, and the *eval92* partition (333 utterances) as the test set. The LibriSpeech corpus, with a total of 960 hours of speech, is used for pre-training only.

The input consists of 40-dimensional log mel filter bank energy (LFBE) features with a window size of 25ms and hop size of 10ms, with per-speaker mean normalization for WSJ but not for LibriSpeech (we do not use any information beyond the audio of LibriSpeech). To speed up training, after data augmentation we stack every 3 consecutive frames.

We investigate the effect of pre-training on phone-based and character-based connectionist temporal classification (CTC) systems [24]. The phone-based system uses a token set of 351 position-dependent phones, generated by the Kaldi *s5* recipe [25]. The character-based system uses 60 characters including the alphabet, digits, and punctuation symbols. Acoustic model training is implemented with TensorFlow [26]; we use its beam search algorithm, with a beam size of 20, for evaluating phone/character error rates.

Our acoustic model consists of 4 bidirectional LSTM layers [27] with 512 units in each direction. For pre-training, the output feature space of $f(X)$ has a dimensionality of 128. The reconstruction network g has two hidden layers of 1024 ReLU [4] units each. We use Adam [28] as the optimizer for both pre-training and fine-tuning, with initial learning rate tuned by grid search, mini-batch size 4 for fine-tuning on *si84* and 16 for *si284*, and maximum number of epochs 50. We apply dropout [5] at all layers, with rate tuned over $\{0.0, 0.1, 0.2, 0.5\}$. We use the development set phone error rate (PER) at the end of each epoch as the criterion for hyperparameter search and early stopping. The learning rate and dropout are tuned once for the supervised baseline, and the resulting values are used in all fine-tuning experiments. For pre-training, optimization parameters are tuned to minimize the dev set reconstruction loss, which happens within 15 epochs. We set the maximum mask widths to $n_F = 8$ and $n_T = 16$, and tune the numbers of masks m_F and m_T based on development set ASR performance.

4.2. Phone-based: Pre-train on *si284*

We first pre-train the acoustic model on *si284* and fine-tune it on *si84*, to investigate the effect of masking parameters used in pre-training. Note that in this setting, there is no domain difference between pre-training and fine-tuning. The supervised baseline yields a dev PER of 18.52%.

Table 1: Dev set %PERs obtained by phone-based systems pre-trained on *si284* and fine-tuned on *si84*, using different numbers of frequency masks (m_F) and time masks (m_T). The baseline PER without pre-training is 18.52%.

| | $m_T = 0$ | $m_T = 1$ | $m_T = 2$ | $m_T = 3$ |
|-----------|-----------|-----------|--------------|-----------|
| $m_F = 0$ | 18.33 | 18.51 | 17.83 | 18.20 |
| $m_F = 1$ | 17.56 | 17.69 | 17.18 | 17.29 |
| $m_F = 2$ | 17.29 | 17.53 | 17.47 | 17.40 |
| $m_F = 3$ | 17.76 | 17.57 | 17.54 | 17.49 |

Table 2: Dev set %PERs of phone-based systems fine-tuned with different amounts of supervised data, and initialized with different pre-trained models.

| | Baseline | Pre-train <i>si284</i> | Pre-train <i>Libri.</i> w/o LIN | Pre-train <i>Libri.</i> w/ LIN |
|--------------|----------|---------------------------|---------------------------------------|--------------------------------------|
| <i>si84</i> | 18.52 | 17.18 | 17.61 | 17.31 |
| + SpecAug | 16.83 | 15.56 | 15.64 | 14.92 |
| <i>si284</i> | 9.16 | 9.23 | 9.15 | 8.50 |
| + SpecAug | 7.98 | 8.21 | 8.19 | 7.46 |

Table 1 gives the dev set PERs after fine-tuning, with $m_F, m_T \in \{0, 1, 2, 3\}$. The case $m_F = m_T = 0$ corresponds to reconstructing all of the input spectrogram, which reduces to the normal auto-encoder objective, and does not significantly improve the acoustic model. This indicates that it is hard for the standard auto-encoder approach to learn useful representations with this bidirectional architecture, perhaps because given the full context, the reconstruction problem becomes too easy. We also observe that it is important to have at least one frequency mask, demonstrating the importance of exploring the joint time-frequency structure. To verify the importance of masking segments rather than individual frames or frequency bins, we pre-train another model where the total numbers of masked frames and frequency bins are the same as those of our method using the best parameters ($m_F = 1, m_T = 2$), but without constraining the masks to be contiguous; this model gives a worse dev PER of 17.61%.

Based on the above results, we fix $m_F = 1$ and $m_T = 2$ for pre-training phone-based systems. For this model pre-trained on *si284*, we fine-tune with different amounts of supervised data, with or without augmenting the training set (using SpecAugment). The dev set PERs are given in Table 2 (second column). We observe that pre-training is clearly helpful when the supervised set is small (i.e., *si84*).

4.3. Phone-based: Pre-train on LibriSpeech

We next explore how the amount and domain of the unlabeled data affect performance, by pre-training on LibriSpeech with 960 hours of speech. For pre-training we use mini-batch size 128, and find that early stopping occurs after 7 epochs. Since there is a domain difference between LibriSpeech and WSJ, we also investigate the effect of domain adaptation for fine-

¹LDC catalog numbers LDC93S6B and LDC94S13B.

Table 3: Dev set %CERs of character-based systems pre-trained on LibriSpeech, and fine-tuned with different amounts of supervised data.

| | Baseline | Pre-train <i>Libri.</i> w/o LIN | Pre-train <i>Libri.</i> w/ LIN |
|--------------|----------|---------------------------------------|--------------------------------------|
| <i>si284</i> | 15.23 | 14.02 | 13.29 |
| + SpecAug | 12.98 | 12.26 | 11.70 |
| <i>si284</i> | 7.01 | 6.90 | 6.48 |
| + SpecAug | 6.29 | 6.19 | 5.61 |

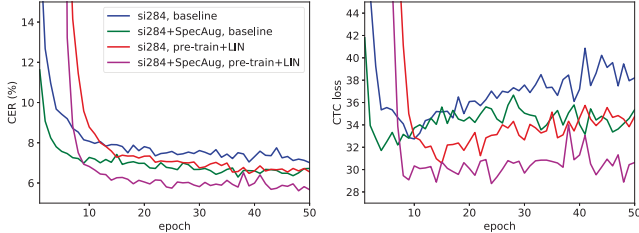


Fig. 2: Dev set learning curves (%CER and CTC loss) of different systems pre-trained on *LibriSpeech*. The first 5 epochs of fine-tuning update only the LIN and softmax layers.

tuning. For domain adaptation we use linear input network (LIN, [29, 30]), which inserts an additional linear layer (initialized as the identity mapping) between the input and the pre-trained network, and only adapts this layer and the softmax layer for the first 5 epochs of supervised training.

The dev set performance when pre-training on LibriSpeech is given in Table 2, with or without LIN adaptation. We observe that without LIN, the performance improvement tends to be smaller than that of pre-training on *si284*. With LIN adaptation, we obtain consistently better PERs, even when fine-tuning on *si284*. Furthermore, the gains from pre-training and SpecAugment are additive.

4.4. Character-based: Pre-train on LibriSpeech

To study how pre-training interacts with different modeling units, we repeat the above experiments for character-based systems. We tune the masking parameters as before (pre-train on *si284* and fine-tune on *si84*), and set $m_F = 3$ and $m_T = 2$ for pre-training on LibriSpeech.

The dev set performance of pre-trained character-based systems is given in Table 3. The observations are consistent with those on the phone-based systems. Fig. 2 shows learning curves of the CTC systems in terms of both CER and the average CTC loss over the dev set, with or without SpecAugment. We see that, although the two criteria (CER and loss) do not synchronize completely, the pre-trained systems are advantageous in terms of both. Note that all four models are trained with the same optimization parameters, and the loss curves with pre-training generally show less overfitting.

Finally, we evaluate word error rates (WERs) for the above character-based systems using the WFST-based frame-

Table 4: %WERs obtained by character-based CTC systems on the test set. Pre-training is done on *LibriSpeech*.

| Method | WER |
|--|-------------|
| EESSEN [31] (extended tri-gram) | 7.34 |
| <i>si284</i> | 7.69 |
| <i>si284</i> + SpecAug | 7.44 |
| <i>si284</i> + pre-train + LIN | 6.66 |
| <i>si284</i> + SpecAug + pre-train + LIN | 6.33 |

work of Miao *et al.* [31], with the extended 4-gram language model built by the Kaldi recipe. After composing the decoding (TLG) graph, we perform beam search using Kaldi’s decode-faster with beam size 20 and acoustic model scale tuned on the dev set. Test set WERs are given in Table 4. For reference, EESSEN’s character-based system obtains a test WER of 7.34% with a different language model, when trained on *si284*. Our results show that the more accurate pre-trained acoustic models also give improved word-level decodings with a language model.

5. CONCLUSIONS

This work demonstrates that pre-training by masked reconstruction leads to consistent performance improvement for CTC-based ASR. Some questions remain open. We have chosen different masking parameters for pre-training the phone-based and character-based systems, by tuning on the development set; it would be good to have a more efficient way of choosing these hyperparameters. In addition, a thorough comparison is needed with other recent work on representation learning approaches [8, 7, 6, 22] to separate the effects of model type versus pre-training approach; this is not trivial as it is not straightforward to extend these prior approaches to bidirectional recurrent models.

6. ACKNOWLEDGEMENTS

The authors would like to thank Yang Chen for helpful discussions, and Amazon for providing computational resources. This material is based upon work supported by the Air Force Office of Scientific Research under award number FA9550-18-1-0166 and by NSF award number 1816627.

7. REFERENCES

- [1] G. E. Hinton, “Training products of experts by minimizing contrastive divergence,” *Neural computation*, vol. 14, no. 8, 2002.
- [2] G. E. Hinton, L. Deng, D. Yu, G. Dahl, A-R. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, T. N. Sainath, and B. Kingsbury, “Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups,” *IEEE Signal Processing Magazine*, 2012.

- [3] X. Glorot and Y. Bengio, “Understanding the difficulty of training deep feedforward neural networks,” in *AISTATS*, 2010.
- [4] V. Nair and G. E. Hinton, “Rectified linear units improve restricted Boltzmann machines,” in *ICML*, 2010.
- [5] N. Srivastava, G. E. Hinton, A. Krizhevsky, I. Sutskever, and R. R. Salakhutdinov, “Dropout: A simple way to prevent neural networks from overfitting,” *JMLR*, 2014.
- [6] S. Pascual, M. Ravanelli, J. Serrà, A. Bonafonte, and Y. Bengio, “Learning problem-agnostic speech representations from multiple self-supervised tasks,” in *INTERSPEECH*, 2019.
- [7] Y.-A. Chung, W.-N. Hsu, H. Tang, and J. Glass, “An unsupervised autoregressive model for speech representation learning,” in *INTERSPEECH*, 2019.
- [8] S. Schneider, A. Baevski, R. Collobert, and M. Auli, “wav2vec: Unsupervised pre-training for speech recognition,” *arXiv preprint arXiv:1904.05862*, 2019.
- [9] M. E. Peters, M. Neumann, M. Iyyer, M. Gardner, C. Clark, K. Lee, and L. Zettlemoyer, “Deep contextualized word representations,” in *NAACL-HLT*, 2018.
- [10] K. Lagler, M. Schindelegger, J. Böhm, H. Krásná, and T. Nilsson, “Gpt2: Empirical slant delay model for radio space geodetic techniques,” *Geophysical research letters*, 2013.
- [11] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” in *NAACL-HLT*, 2019.
- [12] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. R. Bowman, “GLUE: A multi-task benchmark and analysis platform for natural language understanding,” in *ICLR*, 2019.
- [13] D. Park, W. Chan, Y. Zhang, C. Chiu, B. Zoph, E. Cubuk, and Q. V. Le, “SpecAugment: A simple data augmentation method for automatic speech recognition,” *arXiv:1904.08779 [eess.AS]*, Apr. 18 2019.
- [14] E. Dunbar, R. Algayres, J. Karadayi, M. Bernard, J. Benjumea, X.-N. Cao, L. Miskic, C. Dugrain, L. Ondel, A. W. Black, et al., “The zero resource speech challenge 2019: TTS without T,” in *Interspeech*, 2019.
- [15] J. Chorowski, R. J. Weiss, S. Bengio, and A. V. D. Oord, “Unsupervised speech representation learning using wavenet autoencoders,” *TASLP*, vol. 27, 2019.
- [16] H. Kamper, A. Jansen, and S. Goldwater, “A segmental framework for fully-unsupervised large-vocabulary speech recognition,” *Computer Speech & Language*, 2017.
- [17] D. Harwath, A. Torralba, and J. Glass, “Unsupervised learning of spoken language with visual context,” in *NeurIPS*, 2016.
- [18] W.-N. Hsu, Y. Zhang, and J. Glass, “Unsupervised learning of disentangled and interpretable representations from sequential data,” in *NeurIPS*, 2017.
- [19] B. Milde and C. Biemann, “Unspeech: Unsupervised speech context embeddings,” in *Interspeech*, 2018.
- [20] A. V. D. Oord, Y. Li, and O. Vinyals, “Representation learning with contrastive predictive coding,” *arXiv preprint arXiv:1807.03748*, 2018.
- [21] R. D. Hjelm, A. Fedorov, S. Lavoie-Marchildon, K. Grewal, P. Bachman, A. Trischler, and Y. Bengio, “Learning deep representations by mutual information estimation and maximization,” in *ICLR*, 2018.
- [22] Anonymous, “vq-wav2vec: Self-supervised learning of discrete speech representations,” in *Submitted to ICLR*, 2020, under review.
- [23] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, “Librispeech: An asr corpus based on public domain audio books,” in *ICASSP*, April 2015.
- [24] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, “Connectionist temporal classification: Labelling unsegmented sequence data with recurrent neural networks,” in *ICML*, 2006.
- [25] D. Povey, A. Ghoshal, G. Boulianne, L. Burget, O. Glembek, N. Goel, M. Hannemann, P. Motlicek, Y. Qian, P. Schwarz, J. Silovsky, G. Stemmer, and K. Vesely, “The Kaldi speech recognition toolkit,” in *ASRU*, 2011.
- [26] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, and et al., “TensorFlow: Large-scale machine learning on heterogeneous systems,” 2015.
- [27] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural Computation*, 1997.
- [28] D. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *ICLR*, 2015.
- [29] J. Neto, L. Almeida, M. Hochberg, C. Martins, L. Nunes, S. Renals, and T. Robinson, “Speaker-adaptation for hybrid HMM-ANN continuous speech recognition system,” in *EuroSpeech*, 1995.
- [30] K. Yao, D. Yu, F. Seide, H. Su, L. Deng, and Y. Gong, “Adaptation of context-dependent deep neural networks for automatic speech recognition,” in *SLT*, 2012.

- [31] Y. Miao, M. Gawayyed, and F. Metze, “EESSEN: End-to-end speech recognition using deep RNN models and WFST-based decoding,” in *ASRU*, 2015.