# Pretrained Semantic Speech Embeddings for End-to-End Spoken Language Understanding via Cross-Modal Teacher-Student Learning

*Pavel Denisov, Ngoc Thang Vu*

Institute for Natural Language Processing (IMS), University of Stuttgart, Germany

{pavel.denisov, thang.vu}@ims.uni-stuttgart.de

## Abstract

Spoken language understanding is typically based on pipeline architectures including speech recognition and natural language understanding steps. These components are optimized independently to allow usage of available data, but the overall system suffers from error propagation. In this paper, we propose a novel training method that enables pretrained contextual embeddings to process acoustic features. In particular, we extend it with an encoder of pretrained speech recognition systems in order to construct end-to-end spoken language understanding systems. Our proposed method is based on the teacher-student framework across speech and text modalities that aligns the acoustic and the semantic latent spaces. Experimental results in three benchmarks show that our system reaches the performance comparable to the pipeline architecture without using any training data and outperforms it after fine-tuning with ten examples per class on two out of three benchmarks.

**Index Terms**: spoken language understanding, transfer learning, teacher student learning

## 1. Introduction

Recent developments in the fields of electronics, computations and data processing have led to an increased interest in smart assistants with speech interfaces. It is likely driven by the fact that usually people can learn to use speech for interaction intuitively without any special training [1] and make it a primary medium of information exchange. However, speech poses a major challenge to a machine when it comes to the task of extraction of information intended to be transmitted by a human speaker, also known as Spoken Language Understanding (SLU) [2]. The key difficulty here is that speech is highly variable, e.g. depending on room acoustic, and contains rich information about speakers [3]. Some of them are not useful for SLU. The information extraction task is often performed on the text representation using Natural Language Understanding (NLU) methods [4], while Automatic Speech Recognition (ASR) systems [5, 6] convert speech to text. ASR step removes redundant information from the input and provides some kind of normalized form on the output. At the same time it causes loss of potentially useful information that can not be encoded in the text representation, such as prosody, loudness and speech rate. The operation of finding the most probable sequence of words for speech input is computationally expensive. This is partly solved by various heuristics avoiding exploration of less probable hypothesis [7, 8, 9], what in turn introduces additional errors propagated to NLU component. Finally, the sequential design of pipeline approach leads to unavoidable source of latency, because NLU component can not start its work before ASR is finished, and it is not desirable in the interactive context of smart assistant. The problems of pipeline approach described above can be solved by end-to-end SLU methods.

Existing works on end-to-end SLU modeling either focus on supervised downstream tasks, for example dialog act classification [10], intent detection [12], slot filling [14], independent intent detection and domain classification [11] and joint intent detection, domain classification and slot filling [13], or target a generic semantic embedding [15, 16, 17] usually inspired by such successful models as word embeddings Word2Vec [18] and contextual text embeddings BERT [19]. Highly variable and complex nature of speech leads to large amounts of both data and computational resources required for SLU training compared to NLU training, especially for recently popular approach based on contextual embeddings. While data requirements could be satisfied for unsupervised approaches, computational resources are still a problem. Fortunately, most of the modern language processing methods, including ASR and NLU, are based on neural networks and deep learning. Deep learning offers an easy way to transfer knowledge between learned tasks. This technique is referred as transfer learning and it is successfully applied in both ASR [20, 21] and NLU [22, 19]. Therefore, transfer learning should be a promising direction to explore for SLU as well. Several reports [23, 12, 13, 24, 14, 17] indicate that transfer learning from audio modality through pretraining on ASR task or, alternatively, speech autoencoding, is helpful for downstream SLU tasks. Transfer learning from text modality, however, has been applied only for Speech2Vec [16] and SpeechBERT [17] so far.

We propose a novel method that combines parameters transfer from well trained end-to-end ASR systems [25] such as pretrained ESPnet [26] and end-to-end NLU models such as pretrained BERT [19] with Teacher-Student learning [27, 28] for final alignment of SLU output space to NLU output space in order to construct end-to-end SLU model allowing few-shot transfer of downstream tasks from text to speech. By doing so, we enable pretrained end-to-end contextual embeddings such as BERT to process acoustic features. In particular, we aim to generate fixed length vectors with semantic representation from speech segments of variable length. Transfer learning from both text and audio modalities makes our approach mostly similar to [17] and [16]. In this work, we investigate utterance classification task and focus on zero-shot and few-shot cases, but the described method could be adopted to many types of SLU tasks. Although previous works described a number of experiments for such utterance classification tasks as dialog act classification [29] and intent classification [24], and we use the same datasets for the evaluation, we do not compare our results directly to these works, as this is outside of the scope of our work.

## 2. Method

Figure 1 provides the overview of the proposed method. Our SLU model is a combination of two pretrained models. First, we use Encoder block of pretrained end-to-end ASR model
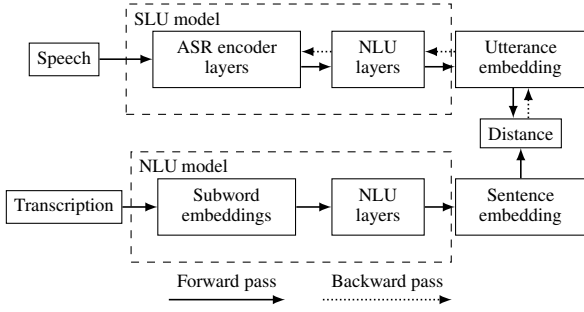
Figure 1: *End-to-end SLU using cross-modal T-S learning.*

[25] in order to covert acoustic features of speech signal to hidden representation. Second, we feed the hidden representation through a learnable linear mapping to pretrained masked language model [19], fine-tuned to produce semantic sentence embedding, which serves as NLU model. Finally, we utilize teacher-student learning method in order to align output of our SLU model to output of pretrained NLU model. Both ASR and NLU models are based on Transformer architecture [30] widely used for sequence processing.

### 2.1. End-to-end ASR

ASR model implements sequence-to-sequence approach and contains two major blocks, encoder and decoder. Encoder takes sequence $X$ with acoustic features and outputs encoded sequence $X_\epsilon$ with hidden representation. Decoder takes the encoded sequence $X_\epsilon$ on input and outputs target sequence $Y$ with text tokens representing transcription of the input utterance. ASR model is trained to minimize weighted sum of cross-entropy objective function calculated from decoder output $\hat{Y}$ and ground truth transcription $Y$ with CTC objective function calculated from learnable linear mapping of encoder output $X_\epsilon$ and ground truth transcription $Y$.

### 2.2. NLU

NLU model is a neural network that takes sequence $X$ with text tokens on input and produces encoded sequence $X_\epsilon$. Semantic sentence embedding vector $y$ is obtained by applying pooling operation to the encoded sequence $X_\epsilon$. The model's parameters are initially pretrained with the tasks of masked token and next sentence prediction from the encoded sequence $X_\epsilon$ representing contextual text token embeddings, as it is done with BERT model [19]. After that, the model is extended with pooling operation over the encoded sequence $X_\epsilon$, producing pooled output $y$, and is fine-tuned on specialized datasets to encode more semantic information to the pooled output $y$.

### 2.3. Teacher-Student learning

Teacher-Student learning minimizes distance-based objective function between outputs of two models on same or equivalent inputs with the aim to update Student model's parameters so that its output becomes more similar to the output of Teacher model. The parameters of Teacher model are not updated during this process. The final stage of our method is the alignment of SLU output to NLU output with Teacher-Student learning method, where SLU model consumes speech recordings and plays the Student role, while NLU model consumes ground truth transcriptions and plays the Teacher role.

## 3. Experimental setup

### 3.1. ASR model

We adopt the latest LibriSpeech recipe [31] from ESPnet toolkit. Transformer network has attention dimension 512, feed-forward inner dimension 2048, 8 heads, 12 blocks in the encoder and 6 blocks in the decoder. Input features are 80-dimensional log Mel filterbank coefficients with 3-dimensional pitch value, frame size is 25 ms and shift is 10 ms. Output labels are 100 subword units, automatically learned with unigram language model algorithm [32] from lowercased concatenation of LibriSpeech and TED-LIUM LM training data with transcriptions of the acoustic training data. Training data combines LibriSpeech, Switchboard, TED-LIUM 3, AMI, WSJ, Common Voice 3, SWC, VoxForge and M-AILABS datasets with a total amount of 3249 hours. Validation data combines validation subsets of LibriSpeech, TED-LIUM 3 and AMI datasets with a total amount of 38 hours. The training is performed on 4 GPUs using Adam optimizer and square root learning rate scheduling [30] with 25,000 warmup steps and learning rate coefficient 10. SpecAugment data augmentation method [33] is applied dynamically during each batch generation. The model is trained for 24 epochs and evaluated on the validation data after each epoch. The final model is obtained by averaging the parameters of the seven best performing models.

### 3.2. NLU model

We use pretrained `bert-base-nli-stsb-mean-tokens` Sentence-BERT model [34]. The model itself is fine-tuned from the well-known pretrained `bert-base-uncased` model [19]. Transformer network has attention dimension 768, feed-forward inner dimension 3072, 12 heads and 12 blocks. Input text is tokenized to 30,000 subword units. The model is pretrained with masked LM and next sentence prediction tasks on BooksCorpus and English Wikipedia datasets. Pooling operation `MEAN` is added to obtain the sentence embedding $y$ from the encoded sequence $X_\epsilon$. The sentence embedding is first fine-tuned on SNLI and MultiNLI datasets for 3-way classification between *contradiction*, *entailment* and *neutral* classes for a given pair of sentences using cross-entropy objective function. After that, the sentence embedding is fine-tuned on STSb dataset for prediction of cosine similarity for a given pair of sentences using mean-squared-error objective function.

### 3.3. SLU model

SLU model is constructed by combining ASR model's encoder with self-attention blocks of NLU model, so that NLU model receives the hidden representation from ASR encoder instead of the output of input embedding layer of NLU model. Linear layer is added between ASR encoder and NLU blocks to map the dimension of hidden representation from 512 to 768. Fine-tuning is performed using Teacher-Student approach by minimizing the distance between output of SLU model for speech recordings and output of NLU model for corresponding transcriptions. We conduct fine-tuning experiments with cosine, L2 and L1 distance based objective functions. SLU model acts as a Student, and we select empirically, which parameters to update during the fine-tuning. NLU model acts as a Teacher, and we freeze its parameters. We employ smaller acoustic dataset consisting of LibriTTS, Common Voice 3, and M-AILABS corpora with a total amount of 1453 hours for the fine-tuning. Our motivation here is to utilize richer transcriptions with punctuation available in these datasets and to supply NLU model with extra

information for potentially semantically finer sentence embeddings. We use the transcriptions as is and do not apply any text preprocessing that is usually done in ASR training, including our end-to-end ASR model. Validation data is the validation subset of LibriTTS corpus with a total duration of 15 hours. We do not apply SpecAugment during the fine-tuning, because it yielded worse results in our early experiments.

### 3.4. Evaluation

SLU model is evaluated on two downstream tasks, dialog act (DA) classification and intent classification, both of which are utterance classification tasks. DA classification is evaluated on two corpora: ICSI Meeting Recorder Dialog Act Corpus (MRDA) and NXT-format Switchboard Corpus (SwDA). Intent classification is evaluated on Fluent Speech Commands (FSC) corpus. Table 1 summarizes the datasets.

Table 1: *SLU evaluation datasets*

| Dataset | Number of classes | Number of utterances | | |
|---|---|---|---|---|
| | | Train | Valid | Test |
| SwBD | 42 | 97,756 | 8,591 | 2,507 |
| MRDA | 6 | 77,596 | 15,721 | 15,398 |
| FSC | 31 | 23,132 | 3,118 | 3,793 |

In order to perform utterance classification, we first train a one layer feed-forward classifier on sentence embeddings, produced by the NLU model from the ground truth transcriptions of training subset, using cross-entropy objective function. After that, we test the classifier on semantic utterance embeddings, extracted from the recordings of testing subset using the SLU model. We report accuracy values as a percentage of correctly classified utterances from the total number of utterances.

### 3.5. Baseline

Traditional approach to SLU tasks is a pipeline of ASR followed by NLU, and we adopt it as a baseline while employing the same ASR and NLU models as in the rest of the experiments. Table 2 reports the results of NLU on ASR output as well as on the ground truth transcriptions. The ground truth results represent an upper bound of accuracy achievable on these datasets with NLU model we use in case of perfect transcriptions on ASR output. The effect of imperfect ASR output varies between datasets depending on the difficulty of recording conditions, the differences between formats of manual transcriptions used to train the classifiers and the tolerance of the downstream tasks to the type of noise that ASR introduces. Amount of errors in ASR output is indicated by Word Error Rate (WER), which is also reported in the table. We select the best performing hyperparameters for the classifier training, but do not fine-tune NLU component for the downstream tasks, because our main goal is SLU as generic speech equivalent for NLU rather then the best possible model for some particular downstream task.

Table 2: *Accuracy of NLU on ASR output and on the ground truth transcriptions and WER of ASR*

| Transcriptions | Accuracy on Test, % | | |
|---|---|---|---|
| | SwBD | MRDA | FSC |
| Ground truth | 71.72 | 77.72 | 100.0 |
| ASR output | 57.23 | 64.06 | 94.57 |
| | WER on Test, % | | |
| ASR output | 28.0 | 29.7 | 7.9 |

## 4. Results

### 4.1. Initial fine-tuning by Teacher-Student learning

#### 4.1.1. Layers for fine-tuning

Our first set of experiments is designed to determine which layers of SLU model should be fine-tuned after the combination of parameters transferred from ASR encoder and NLU. As mentioned before, we insert a linear mapping layer between former ASR and NLU layers because of the difference in dimensionality. It is initialized randomly and its parameters are always updated during the fine-tuning step. In addition to that, we try to fine-tune various amount of layers closest to the mapping layer, meaning top layers of former ASR encoder and bottom layers of former NLU. We do so, because for these layers the output (for ASR encoder) or the input (for NLU) is expected to change after the parameters transfer in contrast to bottom layers of former ASR encoder and top layers of former NLU, where input and output should not change.

We run fine-tuning for 10 epochs using square root learning rate scheduling [30] with 300,000 warmup steps and learning rate coefficient 50, and use cosine distance based objective function. The results are given in Table 3. While it is not completely clear how many layers should be fine-tuned, we can conclusively tell that fine-tuning of former ASR encoder layers is more beneficial than former NLU layers. We decide to fine-tune the two top former ASR encoder layers. The results also illustrate that the optimization of SLU model for smaller distance of its output from the output of NLU model is general enough and translates to accuracy improvements in the downstream tasks, although not in all cases.

Table 3: *Effect of layers fine-tuning*

| ASR layers | NLU layers | Accuracy on Test, % | | | Validation loss |
|---|---|---|---|---|---|
| | | SwBD | MRDA | FSC | |
| 0 | 0 | 43.76 | 56.08 | 68.07 | 0.26 |
| 0 | 1 | 37.61 | 56.47 | 85.53 | 0.19 |
| 1 | 0 | 52.37 | **60.21** | 86.42 | 0.16 |
| 1 | 1 | 52.05 | 58.32 | **86.82** | 0.17 |
| 2 | 0 | 52.93 | 59.42 | 85.76 | **0.15** |
| 3 | 0 | **53.81** | 58.90 | 85.53 | 0.16 |

#### 4.1.2. Learning rate schedule

After deciding which layers to fine-tune, we run a series of experiments to determine the best learning rate schedule. Table 4 presents the combinations of learning rate constant and number of warmup steps explored by us. When we increase number of warmup steps, we notice positive effect from slower learning rate ramp up. However, as number number of warmup steps becomes close to the total number of fine-tuning steps, we have to increase number of epochs from 10 to 20 in order to see the whole fine-tuning process.

Table 4: *Effect of learning rate schedule*

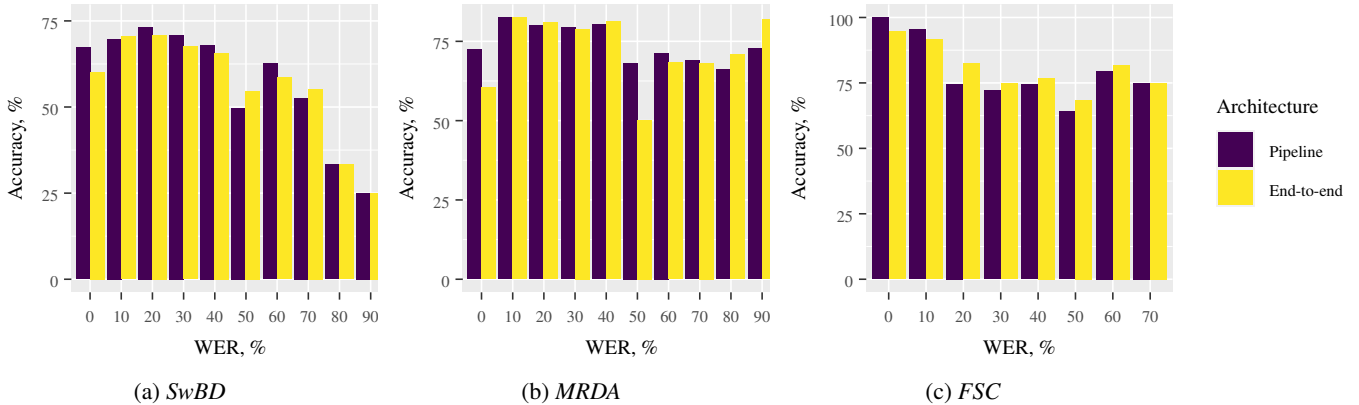| Warmup steps | LR constant | Epochs | Accuracy on Test, % | | | Validation loss |
|---|---|---|---|---|---|---|
| | | | SwBD | MRDA | FSC | |
| 300,000 | 50 | 10 | 52.93 | 59.42 | 85.76 | 0.15 |
| 600,000 | 50 | 10 | 51.18 | 59.95 | 86.84 | 0.14 |
| 600,000 | 50 | 20 | 54.00 | **60.12** | 88.64 | 0.14 |
| 700,000 | 50 | 20 | 51.89 | 58.37 | 88.24 | 0.14 |
| 700,000 | 70 | 20 | 53.73 | 59.67 | 88.08 | 0.14 |
| 700,000 | 30 | 20 | **55.56** | 59.64 | **89.45** | **0.13** |

| (a) *SwBD* | (b) *MRDA* | (c) *FSC* |

Figure 2: *Accuracy comparison for the utterances grouped by ASR WER*

## 4.1.3. Objective function

Comparison of objective functions on downstream tasks, as well as cross-comparison of how selected objective function influences value of others on validation subset, is provided in Table 5. Overall, these results indicate that the evaluated objective functions behave similarly in this task, however L1 distance based objective function yields slightly better results.

Table 5: *Effect of objective function and longer training*

| Objective function | Accuracy on Test, % | | | Validation value | | |
|---|---|---|---|---|---|---|
| | SwBD | MRDA | FSC | Cosine | L2 | L1 |
| Cosine | 55.56 | 59.64 | 89.45 | 0.13 | 0.08 | 0.21 |
| L2 | 53.73 | 59.91 | 88.64 | 0.13 | 0.07 | 0.20 |
| L1 | 56.32 | **60.39** | 89.98 | 0.13 | 0.07 | 0.20 |
| L1, 61 ep. | **58.60** | 60.18 | **91.12** | 0.11 | 0.06 | 0.18 |

## 4.2. Further supervised fine-tuning on downstream tasks

The resulting end-to-end utterance classification model is fully differentiable and can be further optimized for a downstream task by applying standard supervised neural network training methods with few labeled speech samples. This feature should be helpful for the full exploitation of information that is relevant for the task and is encoded in speech, what is less trivial to implement in the traditional ASR and NLU pipeline setup, where intermediate representation has to be discrete and for example flatten the rich variety of prosodic events to few punctuation characters. We examine whether it is useful in practice by running standard supervised classifier training on few samples from training subsets. Table 6 compares the results of fine-tuning of the output layer alone and together with two former ASR encoder layers. We conclude that end-to-end approach indeed can overcome the error propagation problem of pipeline SLU approach by the automatic propagation of the error signal back to relevant parts of SLU system. However, additional training samples may sometimes easily skew the small training dataset away from the testing dataset and cause worse results, so more attention should be payed to the selection of training samples.

## 5. Qualitative Analysis

We attempt to assess the differences between the pipeline and end-to-end SLU approaches in greater detail by looking at the

Table 6: *Effect of supervised fine-tuning on downstream tasks*

| Num. of samples per class | Fine-tuned layers (accuracy on Test, %) | | | | | |
|---|---|---|---|---|---|---|
| | Output layer | | | Output and hidden layers | | |
| | SwBD | MRDA | FSC | SwBD | MRDA | FSC |
| 0 | 58.60 | 60.18 | 91.12 | 58.60 | 60.18 | 91.12 |
| 1 | 58.60 | 60.59 | 93.62 | 58.60 | 60.41 | 94.15 |
| 2 | 58.60 | 60.22 | 93.44 | 58.60 | 60.40 | 95.04 |
| 3 | 58.83 | 60.22 | 93.33 | 58.83 | 60.16 | 94.83 |
| 4 | 58.55 | 60.35 | **93.96** | 58.71 | 60.47 | **95.54** |
| 10 | **60.14** | **60.94** | 93.88 | **60.22** | **61.32** | 95.49 |

accuracy values on groups of utterances split by the WER levels of the baseline ASR system. Our hypothesis is that the pipeline system would make more mistakes on the more challenging recordings characterized by higher WER values, because ASR systems are optimized for phonetic or graphemic similarity to the ground truth and are more likely to lose semantic information in case of errors. Figure 2 shows the accuracy values of the pipeline system and the best end-to-end SLU model (without fine-tuning). The utterances are grouped by WER of the baseline ASR in 10% ranges. The ranges with WER > 100% are not included for brevity, they account for 487, 3847 and 44 utterances in SwBD, MRDA and FSC testing subsets respectively. The results confirm fully our hypothesis on FSC dataset and to some extent on other two datasets.

## 6. Conclusions

We proposed to combine parameters transfer from well trained ASR and NLU models with Teacher-Student learning for final alignment of SLU output space to NLU output space in order to construct end-to-end SLU model allowing few-shot transfer of downstream tasks from text to speech. We outlined necessary steps and settings for the practical pretrained NLU model adaptation in SLU via cross-modal transfer. Our system reaches accuracy of 58.60%, 60.18% and 91.12% on SwBD, MRDA and FSC datasets without fine-tuning and 60.22%, 61.32% and 95.49% after fine-tuning on ten labeled samples per class compared to 57.23%, 64.06% and 94.57% reached by the pipeline system. The results of this research support the idea that text pretrained contextual embeddings can be useful for tasks outside of text modality. The present study also adds new tasks to the growing body of research of language processing methods using Transformer neural networks.

# 7. References

[1] S. Pinker and P. Bloom, "Natural language and natural selection," *Behavioral and brain sciences*, vol. 13, no. 4, pp. 707–727, 1990.

[2] R. De Mori, F. Bechet, D. Hakkani-Tur, M. McTear, G. Riccardi, and G. Tur, "Spoken language understanding," *IEEE Signal Processing Magazine*, vol. 25, no. 3, pp. 50–58, 2008.

[3] J. L. Kröger, O. H.-M. Lutz, and P. Raschke, "Privacy implications of voice and speech analysis–information disclosure by inference," in *IFIP International Summer School on Privacy and Identity Management*. Springer, 2019, pp. 242–258.

[4] C. D. Manning and H. Schütze, *Foundations of statistical natural language processing*. MIT press, 1999.

[5] J. Baker, "The dragon system–an overview," *IEEE Transactions on Acoustics, speech, and signal Processing*, vol. 23, no. 1, pp. 24–29, 1975.

[6] F. Jelinek, *Statistical methods for speech recognition*. MIT press, 1997.

[7] A. Viterbi, "Error bounds for convolutional codes and an asymptotically optimum decoding algorithm," *IEEE transactions on Information Theory*, vol. 13, no. 2, pp. 260–269, 1967.

[8] J. Chorowski and N. Jaitly, "Towards better decoding and language model integration in sequence to sequence models," *Proc. Interspeech 2017*, pp. 523–527, 2017.

[9] T. Hori, J. Cho, and S. Watanabe, "End-to-end speech recognition with word-based rnn language models," in *2018 IEEE Spoken Language Technology Workshop (SLT)*. IEEE, 2018, pp. 389–396.

[10] D. Ortega and N. T. Vu, "Lexico-acoustic neural-based models for dialog act classification," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 6194–6198.

[11] D. Serdyuk, Y. Wang, C. Fuegen, A. Kumar, B. Liu, and Y. Bengio, "Towards end-to-end spoken language understanding," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 5754–5758.

[12] Y.-P. Chen, R. Price, and S. Bangalore, "Spoken language understanding without speech recognition," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 6189–6193.

[13] P. Haghani, A. Narayanan, M. Bacchiani, G. Chuang, N. Gaur, P. Moreno, R. Prabhavalkar, Z. Qu, and A. Waters, "From Audio to Semantics: Approaches to end-to-end spoken language understanding," in *2018 IEEE Spoken Language Technology Workshop (SLT)*. IEEE, 2018, pp. 720–726.

[14] N. Tomashenko, A. Caubrière, and Y. Estève, "Investigating adaptation and transfer learning for end-to-end spoken language understanding from speech," *Proc. Interspeech 2019*, pp. 824–828, 2019.

[15] Y.-A. Chung and J. Glass, "Speech2vec: A sequence-to-sequence framework for learning word embeddings from speech," *arXiv preprint arXiv:1803.08976*, 2018.

[16] Y.-A. Chung, W.-H. Weng, S. Tong, and J. Glass, "Unsupervised cross-modal alignment of speech and text embedding spaces," in *Advances in Neural Information Processing Systems*, 2018, pp. 7354–7364.

[17] Y.-S. Chuang, C.-L. Liu, and H.-Y. Lee, "Speechbert: Cross-modal pre-trained language model for end-to-end spoken question answering," *arXiv preprint arXiv:1910.11559*, 2019.

[18] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," in *Advances in neural information processing systems*, 2013, pp. 3111–3119.

[19] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.

[20] N. T. Vu and T. Schultz, "Multilingual multilayer perceptron for rapid language adaptation between and across language families." in *Interspeech*, 2013, pp. 515–519.

[21] J. Kunze, L. Kirsch, I. Kurenkov, A. Krug, J. Johannsmeier, and S. Stober, "Transfer learning for speech recognition on a budget," in *Proceedings of the 2nd Workshop on Representation Learning for NLP*, 2017, pp. 168–177.

[22] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, "Language models are unsupervised multitask learners," 2019.

[23] Y. Qian, R. Ubale, V. Ramanaryanan, P. Lange, D. Suendermann-Oeft, K. Evanini, and E. Tsuprun, "Exploring asr-free end-to-end modeling to improve spoken language understanding in a cloud-based dialog system," in *2017 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2017, pp. 569–576.

[24] L. Lugosch, M. Ravanelli, P. Ignoto, V. S. Tomar, and Y. Bengio, "Speech Model Pre-training for End-to-End Spoken Language Understanding," *arXiv preprint arXiv:1904.03670*, 2019.

[25] S. Kim, T. Hori, and S. Watanabe, "Joint ctc-attention based end-to-end speech recognition using multi-task learning," in *2017 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, 2017, pp. 4835–4839.

[26] P. Denisov and N. T. Vu, "Ims-speech: A speech to text tool," *Studientexte zur Sprachkommunikation: Elektronische Sprachsignalverarbeitung 2019*, pp. 170–177, 2019.

[27] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," *arXiv preprint arXiv:1503.02531*, 2015.

[28] J. Li, M. L. Seltzer, X. Wang, R. Zhao, and Y. Gong, "Large-scale domain adaptation via teacher-student learning," *Proc. Interspeech 2017*, pp. 2386–2390, 2017.

[29] D. Ortega, C.-Y. Li, G. Vallejo, P. Denisov, and N. T. Vu, "Context-aware neural-based dialog act classification on automatically generated transcriptions," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 7265–7269.

[30] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in neural information processing systems*, 2017, pp. 5998–6008.

[31] S. Karita, N. Chen, T. Hayashi, T. Hori, H. Inaguma, Z. Jiang, M. Someki, N. E. Y. Soplin, R. Yamamoto, X. Wang *et al.*, "A comparative study on transformer vs rnn in speech applications," *arXiv preprint arXiv:1909.06317*, 2019.

[32] T. Kudo, "Subword regularization: Improving neural network translation models with multiple subword candidates," *arXiv preprint arXiv:1804.10959*, 2018.

[33] D. S. Park, W. Chan, Y. Zhang, C.-C. Chiu, B. Zoph, E. D. Cubuk, and Q. V. Le, "SpecAugment: A Simple Data Augmentation Method for Automatic Speech Recognition," *Proc. Interspeech 2019*, pp. 2613–2617, 2019.

[34] N. Reimers and I. Gurevych, "Sentence-bert: Sentence embeddings using siamese bert-networks," in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 11 2019.