

# AIPNET: GENERATIVE ADVERSARIAL PRE-TRAINING OF ACCENT-INVARIANT NETWORKS FOR END-TO-END SPEECH RECOGNITION

Yi-Chen Chen<sup>1,2</sup>, Zhaojun Yang<sup>1</sup>, Ching-Feng Yeh<sup>1</sup>, Mahaveer Jain<sup>1</sup>, Michael L. Seltzer<sup>1</sup>

<sup>1</sup>Facebook AI, USA

<sup>2</sup>Speech Processing and Machine Learning Laboratory, National Taiwan University, Taiwan

## ABSTRACT

As one of the major sources in speech variability, accents have posed a grand challenge to the robustness of speech recognition systems. In this paper, our goal is to build a unified end-to-end speech recognition system that generalizes well across accents. For this purpose, we propose a novel pre-training framework AIPNet based on generative adversarial nets (GAN) for accent-invariant representation learning: **Accent Invariant Pre-training Networks**. We pre-train AIPNet to disentangle accent-invariant and accent-specific characteristics from acoustic features through adversarial training on accented data for which transcriptions are not necessarily available. We further fine-tune AIPNet by connecting the accent-invariant module with an attention-based encoder-decoder model for multi-accent speech recognition. In the experiments, our approach is compared against four baselines including both accent-dependent and accent-independent models. Experimental results on 9 English accents show that the proposed approach outperforms all the baselines by 2.3 ~ 4.5% relative reduction on average WER when transcriptions are available in all accents and by 1.6 ~ 6.1% relative reduction when transcriptions are only available in US accent.

**Index Terms**— Generative adversarial network, end-to-end speech recognition, accent-invariance

## 1. INTRODUCTION

Accents are defined as variations in pronunciation within a language and are often peculiar to geographical regions, individuals, social groups, etc. As one of the major sources in speech variability, accents have posed a grand technical challenge to the robustness of ASR systems. Due to the acoustic discrepancy among accents, an ASR system that is trained on the speech data of one accent (e.g., native) often fails to recognize speech of another unseen accent (e.g., non-native). In this work, we focus on learning accent-invariant representations, aiming to build a universal ASR system that is generalizable across accents.

There is an extensive literature on multi-accent modeling for speech recognition [1] [2]. The existing approaches can be categorized into two classes in general: accent-independent and accent-dependent. Accent-independent modeling focuses on building a universal model that generalizes well across accents. One popular baseline is to train a model on all the data of different accents [3] [4] [5]. Elfeky *et al.* have attempted to build a unified multi-accent recognizer from a pre-defined unified set of CD states by learning from the ensemble of accent-specific models [3]. Yang *et al.* have proposed to jointly model ASR acoustic model and accent identification classifier through multi-task learning [6]. Accent-dependent approaches either take accent-related information, such as accent embedding or i-vectors, as an complementary input in the modeling or

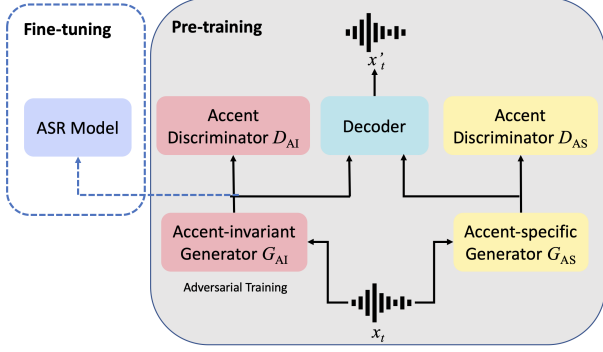
adapt a unified model on accent-specific data [7] [8] [9] [10] [11]. Accent-dependent models usually outperform the unified ones with known accent labels or on an accent-specific dataset, while accent-independent models demonstrate better generalizability on average when accent labels are unavailable during testing.

Generative adversarial nets (GAN) [12] or gradient reverse technique [13] has gained popularity in learning a representation that is invariant to domains or conditions [14] [15] [16] [17]. Serdyuk *et al.* have applied adversarial training to generate noise-invariant representations for speech recognition [15]. Gradient reversal training has recently been used for learning domain-invariant features to alleviate the mismatch between accents during training [16]. Bousmalis *et al.* have proposed to a GAN-based pixel-level transformation from one domain to the other and have shown great improvement over state-of-the-art on a number of domain adaptation tasks [17].

This work focuses on learning accent-invariance with the goal of building a unified accent-independent system for end-to-end speech recognition. Pre-training has shown its superiority in many computer vision and NLP tasks [18] [19] [20], while research efforts on accent model pre-training thus far have been limited. We propose a novel pre-training framework AIPNet based on GAN for accent-invariant representation learning: **Accent Invariant Pre-training Networks**. Unlike most of the existing work that unites the modeling of acoustics and accents in a single stage, our approach decouples accent modeling from acoustic modeling and consists of two stages: pre-training and fine-tuning. In the pre-training stage, AIPNet is built through adversarial training to disentangle accent-invariant and accent-specific characteristics from acoustic features. As transcriptions are not needed in pre-training, AIPNet allows us to make use of many possible accent resources for which transcriptions are unavailable. In the fine-tuning stage, we adopt an attention-based encoder-decoder model for sequence-to-sequence speech recognition. Specifically, we plug in the accent-invariant embeddings in AIPNet into ASR model for further optimization. Experimental results on 9 English accents show significant WER reduction compared to four popular baselines, indicating the effectiveness of AIPNet on accent-invariance modeling. As a general framework for learning domain-invariance, AIPNet can be easily generalized to model any variabilities, such as speakers or speech noise, in addition to accents.

## 2. AIPNET

In this section we describe AIPNet in details. Our approach consists of two stages: pre-training and fine-tuning. In the pre-training stage, the model is built through adversarial training with the goal of learning accent-invariant representations. In the fine-tuning stage, we stack the pre-trained model with downstream tasks for further



**Fig. 1:** The framework of AIPNet including both pre-training and fine-tuning stages.

optimization. In this work, we use end-to-end ASR as a downstream application, focusing on improving accent robustness for speech recognition. The framework of AIPNet is illustrated in Fig. 1.

Suppose the input is an utterance  $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T)$ , where  $\mathbf{x}_t$  represents the feature vector at time step  $t$ . The speaker accent corresponding to  $\mathbf{x}_t$  is denoted as  $a_t \in \{1, 2, \dots, C\}$ , where  $C$  is the number of accents in the training data.

## 2.1. Accent-Invariance Pre-training

The goal of pre-training is to learn accent-invariant representations from accented training data. We define three types of losses for this purpose, including adversarial loss to disentangle accent-invariant and accent-specific information, reconstruction loss to enforce acoustic characteristics to be preserved in the disentangled representations, as well as consistency regularization to detach linguistic information from accent-specific representations.

### 2.1.1. Adversarial Loss

To learn accent-invariant representations, we define two mappings from speech data: accent-invariant generator  $G_{AI}(\mathbf{x}_t)$  and accent-specific generator  $G_{AS}(\mathbf{x}_t)$ . We also define two discriminators  $D_{AI}(G_{AI})$  and  $D_{AS}(G_{AS})$  that output probabilities of accents to ensure that  $G_{AI}$  and  $G_{AS}$  encode the corresponding information. Specifically, we train  $D_{AI}$  and  $D_{AS}$  to maximize the probability of assigning correct accent labels to samples from  $G_{AI}$  and  $G_{AS}$  respectively, *i.e.*, minimizing cross-entropy loss  $L_{CE}^{AI}$  and  $L_{CE}^{AS}$ :

$$\min_{D_{AS}, G_{AS}} L_{CE}^{AS} = \sum_{t=1}^T -\log P(a_t | G_{AS}(\mathbf{x}_t)), \quad (1)$$

$$\min_{D_{AI}} L_{CE}^{AI} = \sum_{t=1}^T -\log P(a_t | G_{AI}(\mathbf{x}_t)). \quad (2)$$

To decouple accent-related information from  $G_{AI}$ , we simultaneously train  $G_{AI}$  such that  $D_{AI}$  is confused about accent labels of samples from  $G_{AI}$ . The objective is to maximize cross-entropy loss  $L_{CE}^{AI}$ , equivalent to minimize the negative cross-entropy:

$$\min_{G_{AI}} -L_{CE}^{AI} = \sum_{t=1}^T \log P(a_t | G_{AI}(\mathbf{x}_t)). \quad (3)$$

### 2.1.2. Reconstruction Loss

The adversarial loss defined between  $D_{AI}$  and  $G_{AI}$  enforces that accent-specific information is disentangled from  $G_{AI}$  but preserved in  $G_{AS}$ . To ensure acoustics characteristics are encoded in the representations from both generators, we further define a decoder with autoencoding structure to reconstruct speech feature  $\mathbf{x}_t$  as  $\mathbf{x}'_t$  from the concatenation of  $G_{AI}(\mathbf{x}_t)$  and  $G_{AS}(\mathbf{x}_t)$ . The decoder is trained by minimizing the reconstruction error  $L_R$ :

$$\min_{decoder, G_{AI}, G_{AS}} L_R = \sum_{t=1}^T \|\mathbf{x}'_t - \mathbf{x}_t\|_2^2. \quad (4)$$

### 2.1.3. Consistency Regularization

Accent-specific attributes are generally stable within an utterance while linguistic-related acoustics have larger intra-utterance variance across time frames. Inspired by the utterance-level stability of accent-specific attributes, we impose a consistency regularization for  $G_{AS}(\mathbf{x}_t)$  such that accent-specific representations from  $G_{AS}$  are consistent across time frames within an utterance:

$$\min_{G_{AS}} L_{CR} = \sum_{t=1}^{T-1} \|G_{AS}(\mathbf{x}_{t+1}) - G_{AS}(\mathbf{x}_t)\|_2^2. \quad (5)$$

This regularization reinforces the preservation of accent-specific information in  $G_{AS}$  meanwhile implicitly encourages linguistic content to be disentangled from  $G_{AS}$ . The multi-scale nature of information in speech data has also been applied in voice conversion and speech denoising [21].

### 2.1.4. Iterative Training

Given the minmax two-player game between  $D_{AI}$  and  $G_{AI}$ , AIPNet pre-training is designed of repeating the following two steps in an iterative manner.

- Update the discriminator  $D_{AI}$  by minimizing  $L_D$ ,
- Freeze the discriminator  $D_{AI}$  and update the rest of the network by minimizing  $L_G$ ,

$$L_D = L_{CE}^{AI}, \quad (6)$$

$$L_G = -L_{CE}^{AI} + \lambda_1 L_{CE}^{AS} + \lambda_2 L_R + \lambda_3 L_{CR}, \quad (7)$$

where  $\lambda$ s are hyper-parameters.

## 2.2. Fine-tuning for End-to-End Speech Recognition

In the fine-tuning stage, the outputs of  $G_{AI}$  which encode accent-invariant linguistic content can be plugged in as inputs of any downstream speech tasks that aim to improve accent robustness, as shown in Fig. 1. In this work, we focus on multi-accent speech recognition and adopt Listen, attend and spell (LAS), a popular attention-based encoder-decoder model [22] for sequence-to-sequence speech recognition. LAS consists of an encoder encoding an input sequence into high-level representations as well as an attention-based decoder generating a sequence of labels from the encoded representations. The encoder is typically a unidirectional or bidirectional LSTM and the decoder is a unidirectional LSTM.

The label inventory for LAS modeling consists of 200 word pieces and is further augmented with two special symbols  $\langle \text{sos} \rangle$  and  $\langle \text{eos} \rangle$  indicating the start of a sentence and the end of a sentence respectively. LAS models the posterior probability of a label

**Table 1:** WER (%) of different approaches in each accent in supervised setting. F1 indicates fine-tuning with  $L_{ASR}$ ; F2 indicates fine-tuning with  $L'_G$ ; AI indicates accent-independent model; AD indicates accent-dependent model.

Approach			Ave.	US	Non-US	CA	FR	IN	KR	PH	LA	GB	VN
Baselines	B1	AI	8.7	5.7	9.0	6.4	8.4	11.2	<b>9.9</b>	7.2	<b>7.8</b>	8.0	13.0
	B2	AI	8.8	<b>5.0</b>	9.1	6.6	9.3	11.0	10.3	<b>6.7</b>	8.1	8.1	12.9
	B3	AD	8.6	5.4	8.9	6.7	8.5	10.9	10.0	6.8	8.6	7.9	<b>12.0</b>
	B4	AI	8.8	5.8	9.1	6.1	8.5	11.7	10.7	7.4	8.4	<b>7.8</b>	<b>12.0</b>
AIPNet	F1	AI	<b>8.4</b>	5.6	<b>8.7</b>	<b>6.0</b>	<b>8.1</b>	<b>9.9</b>	10.3	6.9	8.0	<b>7.8</b>	12.4
	F2	AI	10.1	6.2	10.5	7.9	10.1	12.8	12.1	8.2	9.5	9.4	13.9

**Table 2:** WER (%) of different approaches in each accent in semi-supervised setting. F1 indicates fine-tuning with  $L_{ASR}$ ; F2 indicates fine-tuning with  $L'_G$ ; PL indicates pseudo labeling; AI indicates accent-independent model; AD indicates accent-dependent model.

Approach			Ave.	US	Non-US	CA	FR	IN	KR	PH	LA	GB	VN
w/o PL	Baseline	B1 AI	29.9	10.6	31.8	22.0	33.1	41.0	33.1	28.2	28.7	28.3	40.6
	AIPNet	F1 AI	27.9	<b>9.4</b>	29.8	20.1	30.8	39.0	32.8	25.5	26.4	25.3	39.2
w/ PL	Baselines	B1 AI	26.2	10.3	27.8	18.6	28.3	36.1	29.6	24.8	25.1	24.4	35.8
		B2 AI	25.9	<b>9.4</b>	27.6	19.0	27.7	36.5	29.6	24.2	23.8	25.1	34.9
		B3 AD	25.9	9.6	27.5	19.5	28.0	36.4	29.1	23.7	24.2	24.8	35.0
		B4 AI	25.0	9.7	26.5	<b>18.1</b>	26.7	34.9	28.3	23.7	23.4	23.7	33.6
w/ PL	AIPNet	F1 AI	25.7	12.1	27.0	19.7	27.4	34.7	28.9	23.0	23.6	24.5	34.6
		F2 AI	<b>24.6</b>	11.8	<b>25.9</b>	19.0	<b>26.0</b>	<b>32.6</b>	<b>28.0</b>	<b>22.2</b>	<b>22.8</b>	<b>23.1</b>	<b>33.5</b>

sequence  $\mathbf{y}$  given the input feature sequence  $G_{AI}(\mathbf{X})$  and the previous label history  $\mathbf{y}_{1:j-1}$ :

$$P(\mathbf{y}|G_{AI}(\mathbf{X})) = \prod_{j=1}^J P(y_j|G_{AI}(\mathbf{X}), \mathbf{y}_{1:j-1}). \quad (8)$$

Both encoder and decoder can be trained jointly for speech recognition by maximizing the log probability or minimizing  $L_{ASR}$ :

$$L_{ASR} = \sum_{j=1}^J -\log P(y_j|G_{AI}(\mathbf{X}), \mathbf{y}_{1:j-1}). \quad (9)$$

There are two ways of fine-tuning: 1) fine-tune  $G_{AI}$  and LAS with  $L_{ASR}$ . This requires only transcriptions in the training data; 2) continue with adversarial training as described in Sec 2.1.4 with  $L'_G = L_G + \lambda_4 L_{ASR}$ . This requires both transcriptions and accent labels in the training data. In the experiments, we report results of using both ways.

### 3. EXPERIMENTS

#### 3.1. Dataset

The dataset used in experiments contains utterances in a variety of domains, such as weather or music, collected through crowdsourced workers. There are 9 English accents in total in the dataset, including United States (US), Korea (KR), Philippines (PH), Canada (CA), India (IN), France (FR), Britain (GB), Vietnam (VN) and Latin America (LA). The training set contains 4M (3.8K hours) utterances among which 1% is split as validation data. Particularly, there are 1M and 780K utterances in US and LA respectively and about 330K data in each of the remaining accents. The testing set has 10.8K utterances with 1.2K utterances in each accent. In both training and testing sets, we extract acoustic features using 80-dimensional log Mel-filterbank energies that are computed over a 25ms window every 10ms.

#### 3.2. Experimental Setup

The architecture of each module in AIPNet is a multi-layer LSTM. Specifically, we represent  $G_{AI}$ ,  $G_{AS}$  and decoder using 2 LSTM layers with a hidden size of 768, 256 and 1024 respectively.  $D_{AI}$  and  $D_{AS}$  are represented by a LSTM layer with softmax outputs. The configuration of LAS includes a 4-layer LSTM encoder and a 2-layer LSTM decoder, each with a hidden size of 1024. The hyperparameters  $(\lambda_1, \lambda_2, \lambda_3, \lambda_4)$  are swept within the range  $[0.1, 30]$ . Our experiments have shown that the final results are generally stable across different hyperparameter settings. For simplicity, we report results with  $(\lambda_1, \lambda_2, \lambda_3, \lambda_4) = (1, 10, 10, 10)$  in this paper. We use batch size of 16,000 tokens with 32 GPUs for training. We use Adam with learning rate of  $5 \times 10^{-4}$  in pre-training and  $2.5 \times 10^{-4}$  in fine-tuning,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ . A dropout rate of 0.1 is applied to all the layers. We pre-train AIPNet for 15 epochs and fine-tune LAS for 20 epochs. During inference, speech features are fed into  $G_{AI}$  that is absorbed as part of LAS encoder and outputs of LAS are decoded using beam size of 20 without any external language model.

#### 3.3. Baselines

We compare our approach against four popular baselines B1-B4 for multi-accent speech recognition in the experiments. B1 is an accent-independent model which is trained on the data from all the accents. B2 and B3 have shown strong performance on multi-accent speech recognition in [7]. Specifically, we append accent labels at the end of each label sequence and B2 is trained on the updated sequences from all accents. As accent information is not required at inference, B2 is accent-independent. When training B3 which is accent-dependent, we transform accent 1-hot vector into an embedding through a learnable linear matrix and feed the learned embedding into LAS encoder. During B1-B3 training,  $G_{AI}$  is part of LAS encoder containing 6 LSTM layers (see Section 3.2). B4 is the most similar to our ap-

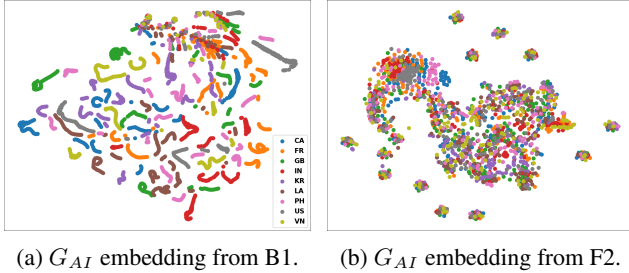


Fig. 2: t-SNE 2-D plots of  $G_{AI}$  embedding from B1 and F2 (w/ PL) in each accent. Each color represents each accent.

proach in spirits, aiming to learn accent-invariant features through gradient reversal [16]. The gradient reversal approach keeps modules of  $G_{AI}$ ,  $D_{AI}$  and ASR model in Fig. 1. Instead of using iterative training in Section 2.1.4, we add a gradient reversal layer between  $G_{AI}$  and  $D_{AI}$  to reverse the backpropagated gradient for  $G_{AI}$  training. For more details about B4, we refer readers to [16].

### 3.4. Experimental Results

As described in Section 2, AIPNet pre-training requires only accent labels in the training data. This approach hence becomes especially useful when there is a large number of accented data without available transcriptions. We design experiments in two settings, *i.e.*, supervised setting where transcriptions are available in all accents and semi-supervised setting where transcriptions are available only in US accent.

#### 3.4.1. Results in Supervised Setting

Table 1 summarizes the results of different approaches in supervised setting. In our approach, we report results of fine-tuning  $G_{AI}$  and LAS with  $L_{ASR}$  using transcriptions (F1), as well as those of fine-tuning the entire network with  $L_G$  using both transcriptions and accent labels (F2). We can see that fine-tuning with  $L_{ASR}$  (F1) outperforms the baselines by 2.3 ~ 4.5% relative reduction on average WER. Compared to all the baselines, F1 has achieved improvement in CA, FR, GB, and especially IN (9.1 ~ 15.3% reduction) but has shown a mediocre performance in each of the remaining accents.

#### 3.4.2. Results in Semi-supervised Setting

In semi-supervised setting where transcriptions are available only in US accent, we compare the performance between B1 and F1. The results are presented in the first two rows of Table 2. As B2, B3, B4 and F2 require the availability of pairs of transcriptions and accent labels for training, the results of these approaches are not available in such scenario. The results have shown that our approach significantly outperforms the baseline model in all accents, achieving 3.4 ~ 11.3% relative WER reduction.

One popular and effective method for semi-supervised learning is to generate target pseudo labels for unlabeled data using an initial model [23]. To achieve better performance, we generate pseudo transcriptions for non-US training data using the US model. As a result, we are able to follow all the experiments in supervised setting in Section 3.4.1. The results with pseudo labeling (PL) are presented in the last six rows of Table 2. By comparing the performance between models with and without pseudo labeling, we can observe that pseudo labeling has shown significant gains for all the

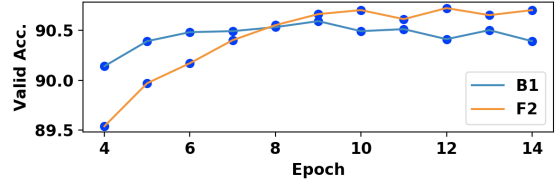


Fig. 3: Word piece validation accuracy of ASR model in B1 and F2.

approaches and almost in each accent, exhibiting its effectiveness on improving generalization performance using unlabeled data. In the scenario with pseudo labeling, fine-tuning with  $L_G$  (F2) outperforms the baselines by 1.6 ~ 6.1% relative reduction on average WER and consistently achieves the best performance in all non-US accents except for CA.

### 3.5. Analysis

In this section, we analyze the properties of AIPNet to better understand its superiority for multi-accent speech recognition. Without loss of generality, we use B1 and F2 (w/ PL) in semi-supervised setting in the analysis.

**Learning accent-invariance** To comprehend the effectiveness of AIPNet on learning accent-invariant representations, we extract embedding (outputs) of  $G_{AI}$  from B1 and F2 respectively for 300 data samples in each accent. Fig. 2 shows t-SNE 2-D visualization of  $G_{AI}$  embedding from B1 (Fig. 2a) and F2 (Fig. 2b) respectively for each accent [24]. As can be seen,  $G_{AI}$  outputs from the baseline B1 tend to be clustered in each accent while those from our approach F2 are mixed across different accents. The visualization demonstrates the validity of the accent-invariant features learned through AIPNet and further explains the better generalization performance that our approach has achieved across accents.

**Reducing overfitting** We further investigate the trend of word piece validation accuracy of the ASR model in B1 and F2, as shown in Fig. 3. Compared to B1, F2 learns more slowly and reaches a better local optimal. The learning objective of F2 consists of both  $L_{ASR}$  and accent-related regularizers (see Section 2). This observation corroborates the effectiveness of the regularization in our approach on reducing the risk of overfitting. It is worth noting that such benefit from the accent-related regularization in fine-tuning is not observed in supervised setting (see Table 1). One possible reason could be that the sufficient labeled training data in supervised setting empowers the ASR model a strong learning capability that might be even weakened by additional regularizations.

## 4. CONCLUSION

In this paper, we proposed AIPNet, a GAN-based pre-training network, for learning accent-invariant representations, aiming to build a unified speech recognition system that generalizes well across accents. As transcriptions are not needed in pre-training, AIPNet provides the flexibility of making use of many possible accent resources for which transcriptions are unavailable. Experiments have shown promising results on 9 English accents compared to the baselines, especially in the case when transcriptions are not available in all accents. Experimental results have demonstrated the effectiveness of AIPNet on learning accent-invariance.

## 5. REFERENCES

- [1] Kanishka Rao and Haşim Sak, “Multi-accent speech recognition with hierarchical grapheme based models,” in *ICASSP*. IEEE, 2017, pp. 4815–4819.
- [2] Hui Lin, Li Deng, Dong Yu, Yi-fan Gong, Alex Acero, and Chin-Hui Lee, “A study on multilingual acoustic modeling for large vocabulary asr,” in *ICASSP*. IEEE, 2009, pp. 4333–4336.
- [3] Mohamed Elfeky, Meysam Bastani, Xavier Velez, Pedro Moreno, and Austin Waters, “Towards acoustic model unification across dialects,” in *SLT*. IEEE, 2016, pp. 624–628.
- [4] Herman Kamper and Thomas Niesler, “Multi-accent speech recognition of afrikaans, black and white varieties of south african english,” in *Interspeech*, 2011.
- [5] Dimitra Vergyri, Lori Lamel, and Jean-Luc Gauvain, “Automatic speech recognition of multiple accented english data,” in *Interspeech*, 2010.
- [6] Xuesong Yang, Kartik Audhkhasi, Andrew Rosenberg, Samuel Thomas, Bhuvana Ramabhadran, and Mark Hasegawa-Johnson, “Joint modeling of accents and acoustics for multi-accent speech recognition,” in *ICASSP*. IEEE, 2018, pp. 1–5.
- [7] Bo Li, Tara N Sainath, Khe Chai Sim, Michiel Bacchiani, Eugene Weinstein, Patrick Nguyen, Zhifeng Chen, Yanghui Wu, and Kanishka Rao, “Multi-dialect speech recognition with a single sequence-to-sequence model,” in *ICASSP*. IEEE, 2018, pp. 4749–4753.
- [8] Mingming Chen, Zhanlei Yang, Jizhong Liang, Yanpeng Li, and Wenju Liu, “Improving deep neural networks based multi-accent mandarin speech recognition using i-vectors and accent-specific top layer,” in *Interspeech*, 2015.
- [9] Sanghyun Yoo, Inchul Song, and Yoshua Bengio, “A highly adaptive acoustic model for accurate multi-dialect speech recognition,” in *ICASSP*. IEEE, 2019, pp. 5716–5720.
- [10] Ngoc Thang Vu, David Imseng, Daniel Povey, Petr Motlicek, Tanja Schultz, and Hervé Bourlard, “Multilingual deep neural network based acoustic modeling for rapid language adaptation,” in *ICASSP*. IEEE, 2014, pp. 7639–7643.
- [11] Yan Huang, Dong Yu, Chaojun Liu, and Yifan Gong, “Multi-accent deep neural network acoustic model with accent-specific top layer using the kld-regularized model adaptation,” in *Interspeech*, 2014.
- [12] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio, “Generative adversarial nets,” in *Advances in neural information processing systems*, 2014, pp. 2672–2680.
- [13] Yaroslav Ganin and Victor Lempitsky, “Unsupervised domain adaptation by backpropagation,” *arXiv preprint arXiv:1409.7495*, 2014.
- [14] Yi-Chen Chen, Sung-Feng Huang, Chia-Hao Shen, Hung-yi Lee, and Lin-shan Lee, “Phonetic-and-semantic embedding of spoken words with applications in spoken content retrieval,” in *SLT*. IEEE, 2018, pp. 941–948.
- [15] Dmitriy Serdyuk, Kartik Audhkhasi, Philémon Brakel, Bhuvana Ramabhadran, Samuel Thomas, and Yoshua Bengio, “Invariant representations for noisy speech recognition,” *arXiv preprint arXiv:1612.01928*, 2016.
- [16] Sining Sun, Ching-Feng Yeh, Mei-Yuh Hwang, Mari Ostendorf, and Lei Xie, “Domain adversarial training for accented speech recognition,” in *ICASSP*. IEEE, 2018, pp. 4854–4858.
- [17] Konstantinos Bousmalis, Nathan Silberman, David Dohan, Dumitru Erhan, and Dilip Krishnan, “Unsupervised pixel-level domain adaptation with generative adversarial networks,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 3722–3731.
- [18] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [19] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” *arXiv preprint arXiv:1810.04805*, 2018.
- [20] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean, “Distributed representations of words and phrases and their compositionality,” in *Advances in neural information processing systems*, 2013, pp. 3111–3119.
- [21] Wei-Ning Hsu, Yu Zhang, and James Glass, “Unsupervised learning of disentangled and interpretable representations from sequential data,” in *Advances in neural information processing systems*, 2017, pp. 1878–1889.
- [22] William Chan, Navdeep Jaitly, Quoc Le, and Oriol Vinyals, “Listen, attend and spell: A neural network for large vocabulary conversational speech recognition,” in *ICASSP*. IEEE, 2016, pp. 4960–4964.
- [23] Dong-Hyun Lee, “Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks,” in *Workshop on Challenges in Representation Learning, ICML*, 2013, vol. 3, p. 2.
- [24] Laurens van der Maaten and Geoffrey Hinton, “Visualizing data using t-sne,” *Journal of machine learning research*, vol. 9, no. Nov, pp. 2579–2605, 2008.