

SPADE: SELF-SUPERVISED PRETRAINING FOR ACOUSTIC DISENTANGLEMENT

John Harvill^{1*}, Jarred Barber^{2†}, Arun Nair², Ramin Pishehvar²

¹University of Illinois Urbana-Champaign, USA

²Amazon Alexa AI, USA

ABSTRACT

Self-supervised representation learning approaches have grown in popularity due to the ability to train models on large amounts of unlabeled data and have demonstrated success in diverse fields such as natural language processing, computer vision, and speech. Previous self-supervised work in the speech domain has disentangled multiple attributes of speech such as linguistic content, speaker identity, and rhythm. In this work, we introduce a self-supervised approach to disentangle room acoustics from speech and use the acoustic representation on the downstream task of device arbitration. Our results demonstrate that our proposed approach significantly improves performance over a baseline when labeled training data is scarce, indicating that our pretraining scheme learns to encode room acoustic information while remaining invariant to other attributes of the speech signal.

Index Terms— keyword spotting, source localization, self-supervised pretraining, disentanglement, acoustics

1. INTRODUCTION

Disentanglement of speech into its multiple components is a fundamental problem in signal processing with applications in voice conversion [1, 2, 3, 4], automatic speech recognition [5, 6, 7, 8], speaker recognition [9], and privacy preservation in speech [10]. The goal of disentanglement is to separate different attributes of the speech signal such that the final signal representation will be invariant to attributes not relevant to the downstream task. Previous work has focused on invariance towards acoustic content¹ so that target attributes like speaker identity or linguistic content will be emphasized. In this paper, we specifically want to preserve acoustic content and extract representations that are invariant to all other attributes in a self-supervised fashion. These representations are then used for the task of device arbitration [11].

The device arbitration task has arisen recently due to the ubiquity of smart voice assistant-enabled devices, which we refer to as “voice assistants” (VA) or “devices” interchangeably. Many households now have multiple VAs in the same room, leading to ambiguity with respect to which device should interact with the user. When the user wishes to begin interaction with a VA, they must first utter a wakeword (“Alexa”, “Hey Google”, etc.) which is then recorded by all devices in the room. For the most natural user experience, only the intended device from the user should wake up and continue to interact with the user. This leads to the device arbitration problem:

given N recordings of a source audio, where each recording comes from a VA, determine which VA is the intended one.

Note that device arbitration is closely related to the well-studied source localization problem [12]. Time Difference Of Arrival (TDOA) [13] is an effective technique that solves source localization for audio signals but unfortunately requires large arrays of microphones not present on VAs. There are several other constraints imposed by modern VAs that make device arbitration a non-trivial problem: (1) The positions of each VA are unknown and could change over time (moving the device from the living room to the kitchen). (2) The acoustic environment of the VAs is unknown. (3) Clock synchronization between devices is not always available, making TDOA between devices infeasible. Given these constraints, the device arbitration task must be solved by only relying on the room acoustic information present in each audio recording.

Motivated by these constraints, we propose SPADE: Self-supervised Pretraining for Acoustic DisEntanglement. SPADE is a pretraining technique that disentangles acoustic information from speech signals by using multiple views of a source audio without the need for labels. We find that when used in combination with previous work on device arbitration [11], SPADE leads to improved performance when less labeled training data is present. Given that labeled data is difficult to collect for this task, SPADE is an invaluable technique for improving performance on device arbitration at zero additional inference cost.

The remainder of our paper is organized as follows: In Section 2 we discuss prior work on device arbitration and related tasks. In Section 3, we discuss the data used in our experiments and the simulation process from which it is generated. In Section 4, we detail our proposed pretraining approach. In Section 5, we discuss our experimental setup and in Section 6 we discuss results. The paper ends with conclusions and suggestions for future work in Section 7.

2. PRIOR WORK

While device arbitration is a relatively new problem, it is closely related to source localization [12, 13]. The goal of source localization is to determine the position of the object emitting sound, i.e. the source. Current techniques rely on large arrays of microphones which are not available given our problem setup. Previous work [11] also demonstrated that directly predicting source distances from each device and making arbitration decisions based on those distances resulted in worse performance.

The goal in arbitration is to select an attended device the user is speaking to based on distance and direction. Common techniques compare signal attributes like Signal-to-Noise Ratio (SNR), estimated distance between source and microphone, cross-correlation, etc. [14]. The goal of device arbitration is different from channel selection [14, 15] since the optimal device is simply defined as that which is closest to the user or as the attended device (defined as

*Work done as part of an internship at Amazon Alexa AI.
† Now at Google Research

¹In this paper, “acoustic content”, or “acoustics” refers to the information in an audio signal related to the Room Impulse Response (RIR) at a particular location in a room when the source audio is played from a different, fixed location.

Parameter	Distribution
Room length/width (m)	Uniform(3.0, 10.0)
Room height (m)	Uniform(2.5, 6.0)
Reverberation time (s)	Beta(1.1, 3.0)
Number of devices	ShiftedPoisson(m=3,l=2,h=15)
Number of noise sources	ShiftedPoisson(m=2,l=1,h=5)
Speech level (dB SPL)	Uniform(55.0, 70.0)
Noise level (dB SPL)	Uniform(50.0, 70.0)

Table 1. Hyperparameters for scene sampling. For the ShiftedPoisson distribution we denote mean with “m”, low with “l” and high with “h.”

the one in the look direction of the user or any other acoustical or visual relevant cue), and not that which leads to better performance on another downstream task. The differences in problem setup and objective between device arbitration and source localization show that device arbitration should be studied separately and has its own set of unique challenges that make it an interesting problem.

We used as a baseline, the machine learning-based device arbitration proposed by Barber et al. [11]. The authors proposed an end-to-end approach to train a neural network to perform device arbitration. Their model consists of a small convolutional feature extractor that runs locally on each device, and a larger arbitration network that runs in the cloud to make the final decision. Results demonstrated significantly improved performance across a variety of room conditions over a simple energy-based approach.

3. DATASET SIMULATION

Currently, there is no large-scale dataset for device arbitration with known ground truth labels, so we run our experiments with simulated data following the three main steps described in Barber et al. [11]:

- **Sample scenario:** Sample a room from a variety of acoustic settings (room length/width, noise sources) as well as device/speaker positions within the room. Given device/speaker positions, generate the arbitration label based on smallest Euclidean distance from device to speaker.
- **Generate RIR:** Given the sampled scenario, generate a Room Impulse Response (RIR) for each device in the scenario using an acoustic simulator.
- **Generate audio:** Convolve speech utterances with generated RIRs and mix with noise for each device. After this step, we have an artificial dataset of device arbitration audio and corresponding ground truth labels.

We use the Image Source Method (ISM) [16] for data simulation and source audio from [11], but have updated the scene sampling hyperparameters to those in Table 1.

4. METHOD

Our arbitration model is based on that proposed by Barber et al. [11] and is composed of two components: the feature encoder and the classifier that makes the arbitration decision. Prior to end-to-end training of the encoder and classifier, we pretrain the encoder using two schemes: contrastive and reconstructive pretraining. Our preprocessing pipeline and pretraining schemes are discussed in the following subsections.

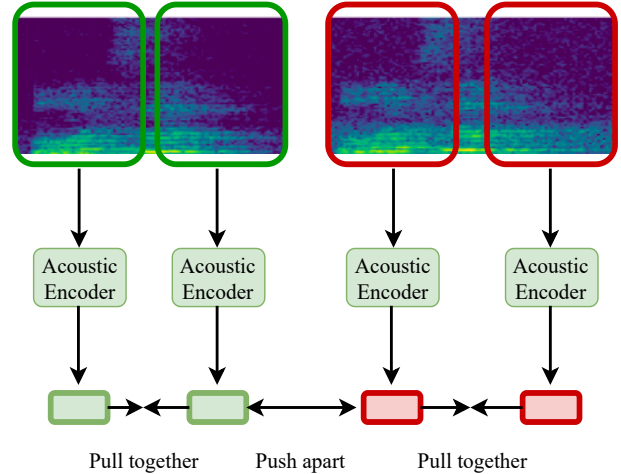


Fig. 1. Contrastive pretraining scheme: Acoustic encoder has shared weights.

4.1. Audio Preprocessing

Audio is first transformed to log-filterbank energy (LFBE) features, where the spectrogram is computed using a 25ms frame size and a frame skip of 10ms. The Mel transform is applied to the spectrogram with 64 Mel bands, followed by the log transform. LFBE features are mean and variance normalized before downstream processing by neural models.

4.2. Encoder Architecture

The encoder architecture is a residual convolutional model composed of 18 convolutional layers, with batch normalization and the ReLU activation function applied to each layer. The encoder model is the same for all approaches discussed in this paper (pretraining, baseline). This network produces a sequential feature representation of the input with much smaller temporal resolution than the input LFBE features. During pretraining we use a small Transformer [17] network to map the sequence of vectors output by the encoder to a single vector. This network is discarded after pretraining but is implicit as part of the encoder in the following pretraining discussions and Figs. 1 and 2.

4.3. Contrastive Pretraining

Contrastive loss functions have been shown to create high-quality representations across a variety of domains like speech and natural language processing [18, 19]. This family of loss functions operates by assuming data are similar or dissimilar with respect to a particular attribute and encouraging embeddings of the data to reflect these relationships via distance in an embedding space. In a device arbitration scene, audio from the same device has the same room acoustic properties, while audio from different devices will have different acoustic properties.² Since we want to encode room acoustic information, we can encourage embeddings of audio from the same device to be similar and simultaneously encourage embeddings of audio from different devices to be orthogonal (see Fig. 1). Since speech content will be the same across all N recordings from a given

²This is due to the assumption that the devices are in different locations.

room, the embedding should be invariant to speech content and only encode acoustic information.

Let us denote the audio recorded by device a as x^a . For this pretraining approach, we create a fixed-length embedding z^a of x^a by passing through the encoder³, i.e. $z^a = \text{encoder}(\text{LFBE}(x^a))$. Then z^a is normalized to have magnitude one. We can denote a continuous slice of x^a from time t_1 to t_2 as $x_{t_1:t_2}^a$ where $t_2 > t_1$ and its corresponding fixed-length embedding as $z_{t_1:t_2}^a$. Given N recordings x^1, x^2, \dots, x^N that are zero-padded to the same length T and a splitting index t_{split} such that $0 < \frac{T}{2} - \epsilon < t_{split} < \frac{T}{2} + \epsilon < T$, where ϵ is a small random jitter, we arrive at our loss function \mathcal{L}_C for one arbitration scenario:

$$\mathcal{L}_1 = \sum_{i=1}^N \sum_{j=1}^N |\langle z_{0:t_{split}}^i, z_{t_{split}:T}^j \rangle - \delta_{ij}| \quad (1)$$

$$\mathcal{L}_2 = \sum_{i=1}^N \sum_{j \neq i} |\langle z_{0:t_{split}}^i, z_{0:t_{split}}^j \rangle| + |\langle z_{t_{split}:T}^i, z_{t_{split}:T}^j \rangle| \quad (2)$$

$$\mathcal{L}_C = \mathcal{L}_1 + \mathcal{L}_2 \quad (3)$$

where δ_{ij} is the Kronecker delta function. Note that \mathcal{L}_1 encourages the two halves of the same audio recording to map to the same embedding. The \mathcal{L}_2 term provides stronger supervision for invariance to speech content by encouraging different recordings of the same respective halves of the audio to be orthogonal.

4.4. Reconstructive Pretraining

Disentanglement of different attributes in speech has been accomplished previously using autoencoding with an information bottleneck [1, 2]. We take a similar approach to disentangle room acoustic information from speech in a self-supervised fashion, making the assumption that the room acoustic properties are constant⁴ for the duration of the wakeword audio (~ 2 s).

Our model consists of a speech encoder, an acoustic encoder, and a reconstruction decoder. The speech encoder $S(\cdot)$ and reconstruction decoder $R(\cdot)$ are Transformers [17] and each produce a sequence of vectors. We design the acoustic encoder $A(\cdot)$ to produce a fixed-dimensional embedding due to the stationarity assumption of the acoustics and the need to create an information bottleneck. Given that N recordings of a source audio all contain identical speech content, the speech representation $s_i = S(x^i)$ should be the same for each recording, i.e. $s^1 \approx s^2 \approx \dots \approx s^N$. To encourage this, we reconstruct one audio recording's LFBE features using its acoustic embedding and the speech embedding from another audio recording in the room. The information bottleneck enforced by creating a small fixed-dimensional embedding encourages the acoustic encoder only to represent information not common to all signals, i.e. room acoustics. Given that the N recordings are not time-aligned, we provide alignment information to the decoder by extracting the envelope $env(x^i)$ of the LFBE features for x^i , which is the mean of the feature vector at each timestep:

$$env(x^i)_k = \frac{1}{N} \sum_{q=1}^N \text{LFBE}(x^i)_{kq} \quad (4)$$

³The Transformer network discussed in Section 4.2 is implied here to follow the encoder to create the fixed-length representation.

⁴This stationarity assumption may not be true in all cases (people/pets moving around), but it is reasonable given that audio is only recorded over a two-second interval.

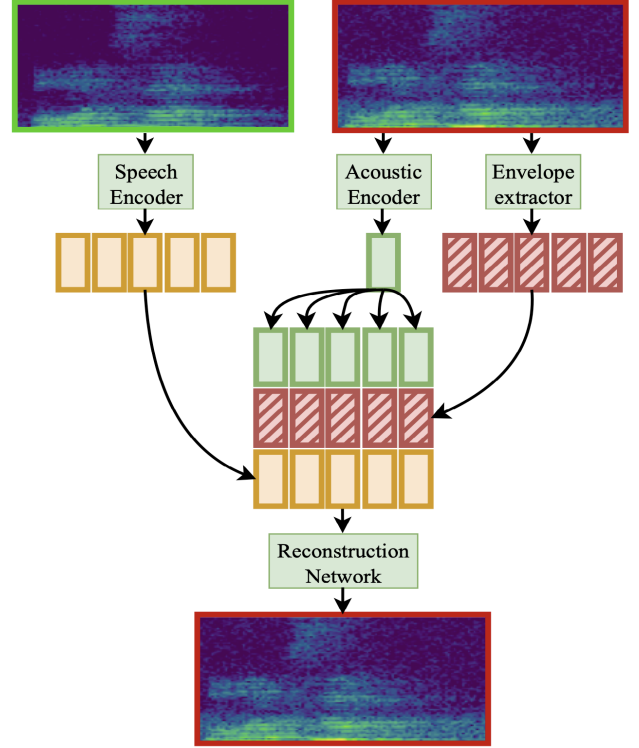


Fig. 2. Reconstruction pretraining scheme: Speech encoder, envelope, and reconstruction network have same time axis dimension as the LFBE time axis. Acoustic encoder creates fixed-length embedding that is copied along time axis of speech encoder and envelope output. Diagonally-shaded units are from non-learnable modules.

where k denotes the time axis, q denotes the feature axis of the LFBE feature matrix, and $N = 64$. Note that the envelope extractor is not a learnable module but rather a simple mean operation at each timestep to produce a low-dimensional representation of the LFBE features (feature vector \rightarrow scalar). We can formally write our loss function \mathcal{L}_R for one arbitration scenario as:

$$\mathcal{L}_R = \frac{1}{N} \sum_{i=1}^N \|\text{LFBE}(x^i) - R(S(x^{\neq i}), A(x^i), env(x^i))\|_2^2 \quad (5)$$

where $x^{\neq i}$ denotes a randomly-selected audio recording from $\{x^1, x^2, \dots, x^N\}$ other than x^i . We form the input to $R(\cdot)$ by copying $A(x^i)$ along the time axis and concatenating it to $S(x^{\neq i})$ and $env(x^i)$. See Fig. 2 for a diagram detailing this process.

4.5. Arbitration Classifier Architecture

The arbitration classifier was implemented previously [11] as a Multilayer Perceptron (MLP) network, but we found further improvement using a self-attention network like the Transformer [17]. For each device i , the encoder outputs a sequence of hidden states $h_1^i, h_2^i, \dots, h_K^i$. For N devices, the hidden states are concatenated along the time axis to form the sequence:

$$H = h_1^1, h_2^1, \dots, h_K^1, h_1^2, h_2^2, \dots, h_K^2, \dots, h_1^N, h_2^N, \dots, h_K^N \quad (6)$$

which is then passed through a network of self-attention layers⁵ to produce the sequence:

$$G = g_1^1, g_2^1, \dots, g_K^1, g_1^2, g_2^2, \dots, g_K^2, \dots, g_1^N, g_2^N, \dots, g_K^N \quad (7)$$

Each sequence $g_1^i, g_2^i, \dots, g_K^i$ is then passed through a second network of self-attention to create a summary G_i over time. Each G_i is then passed through a two-layer feedforward neural network, outputting a scalar logit for device i . The logits are then passed through a softmax layer to produce the arbitration probabilities. The entire classification network is optimized using the crossentropy loss between the arbitration probabilities and the ground truth label distribution.

4.6. Baseline

The contributions of this paper are the self-supervised pretraining approaches, so as a baseline, we use the encoder and classifier networks discussed previously but do not pretrain the encoder.

5. EXPERIMENTS

The training dataset consists of 300k arbitration scenarios. To demonstrate the effectiveness of pretraining for learning representations useful for device arbitration, we create datasets of exponentially decreasing size. For dataset i the training set size $s_i = \lfloor S/4^i \rfloor$, where $S = 300k$ (full training set size). We choose $i \in \{0, 1, 2, 3\}$ such that the smallest training set consists of $\sim 4.7k$ scenarios.

5.1. Experimental Procedure

For each experiment, we choose the final arbitration model based on the checkpoint with the lowest validation loss and evaluate on a held-out test set. We have four experiment setups:

- **Baseline:** Train encoder-classifier model end-to-end on each of the training data subsets of size s_i for $i \in \{0, 1, 2, 3\}$.
- **Contrastive:** Pretrain (acoustic) encoder using contrastive approach on all available training data (300k scenarios, no labels involved). Pick best validation checkpoint as initialization for encoder and then finetune encoder-classifier model end-to-end on each training data subset of size s_i for $i \in \{0, 1, 2, 3\}$.
- **Reconstructive:** Same as contrastive setup except that we use reconstructive pretraining.
- **Combo:** Pretrain using both contrastive and reconstructive pretraining. Loss function becomes $\mathcal{L} = \lambda \mathcal{L}_R + (1 - \lambda) \mathcal{L}_C$ where we set $\lambda = 0.5$. Finetune encoder-classifier model as in previous setups.

6. RESULTS

Results are presented here as *relative error rate* with respect to the performance of the 4.7k baseline setting. Denoting the accuracy of the target method m as acc_m and the accuracy of the 4.7k baseline setting as acc_{base} , we compute relative error rate as:

$$\text{err}_{\text{rel}} = \frac{1 - \text{acc}_m}{1 - \text{acc}_{\text{base}}} \quad (8)$$

⁵Positional encodings are added to the input.

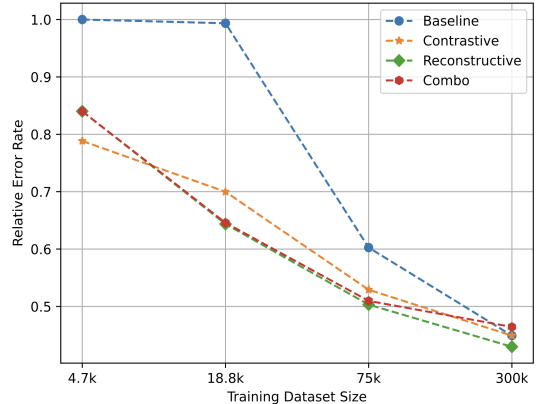


Fig. 3. Relative Error Rate with respect to worst case.

Our results are presented in Figure 3. The most important trend is that the pretraining approaches outperform the baseline by a larger margin when the training dataset is small. This demonstrates that our proposed pretraining schemes preserve acoustic information, learning features relevant to the device arbitration problem and are most beneficial when labeled training data is scarce. Initialization of the encoder with a pretrained checkpoint helps combat overfitting during the device arbitration training process. We also find that the combination of contrastive and reconstructive pretraining does not lead to any noticeable improvement over either approach in isolation, indicating that both approaches may encode similar information.

7. CONCLUSION

In this paper we propose contrastive and reconstructive pretraining, two forms of self-supervised representation learning, that disentangle acoustic content from speech. Unlike previous work that has aimed to create representations that are invariant to acoustic content, we aim to encode acoustic content only and demonstrate its usefulness through the device arbitration problem. We find that both of our proposed pretraining approaches lead to improvement over the baseline, and that improvement is more significant when the labeled training dataset is small. This provides empirical evidence that our pretraining objectives lead to representations of acoustic content that can be useful for the device arbitration task even in the absence of a large training corpus.

Given that self-supervised techniques require no human annotations, it may be possible to apply our proposed approaches to other research problems. For example, our disentangled acoustic representations may be used for other acoustic tasks like room acoustic property estimation or acoustic adaptation for home theater. While not studied in this paper, the speech content representations learned from reconstructive pretraining may be valuable for acoustic-invariant applications like speaker recognition or ASR since they are designed to encode all information except acoustics.

8. REFERENCES

- [1] Kaizhi Qian, Yang Zhang, Shiyu Chang, Xuesong Yang, and Mark Hasegawa-Johnson, "Autovc: Zero-shot voice style transfer with only autoencoder loss," in *International Conference on Machine Learning*. PMLR, 2019, pp. 5210–5219.
- [2] Kaizhi Qian, Yang Zhang, Shiyu Chang, Mark Hasegawa-Johnson, and David Cox, "Unsupervised speech decomposition via triple information bottleneck," in *International Conference on Machine Learning*. PMLR, 2020, pp. 7836–7846.
- [3] Jie Wang, Jingbei Li, Xintao Zhao, Zhiyong Wu, Shiyin Kang, and Helen Meng, "Adversarially learning disentangled speech representations for robust multi-factor voice conversion," *arXiv preprint arXiv:2102.00184*, 2021.
- [4] Jiachen Lian, Chunlei Zhang, and Dong Yu, "Robust disentangled variational speech representation learning for zero-shot voice conversion," in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 6572–6576.
- [5] Kaizhi Qian, Yang Zhang, Heting Gao, Junrui Ni, Cheng-I Lai, David Cox, Mark Hasegawa-Johnson, and Shiyu Chang, "Contentvec: An improved self-supervised speech representation by disentangling speakers," in *International Conference on Machine Learning*. PMLR, 2022, pp. 18003–18017.
- [6] Janek Ebberts, Michael Kuhlmann, Tobias Cord-Landwehr, and Reinhold Haeb-Umbach, "Contrastive predictive coding supported factorized variational autoencoder for unsupervised learning of disentangled speech representations," in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 3860–3864.
- [7] Wei-Ning Hsu, Hao Tang, and James Glass, "Unsupervised adaptation with interpretable disentangled representations for distant conversational speech recognition," *Interspeech*, 2018.
- [8] Sameer Khurana, Shafiq Rayhan Joty, Ahmed Ali, and James Glass, "A factorial deep markov model for unsupervised disentangled representation learning from speech," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 6540–6544.
- [9] Arsha Nagrani, Joon Son Chung, Samuel Albanie, and Andrew Senior, "Disentangled speech embeddings using cross-modal self-supervision," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 6829–6833.
- [10] Ranya Aloufi, Hamed Haddadi, and David Boyle, "Privacy-preserving voice analysis via disentangled representations," in *Proceedings of the 2020 ACM SIGSAC Conference on Cloud Computing Security Workshop*, 2020, pp. 1–14.
- [11] Jarred Barber, Yifeng Fan, and Tao Zhang, "End-to-end alexa device arbitration," in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 926–930.
- [12] Joe C Chen, Kung Yao, and Ralph E Hudson, "Source localization and beamforming," *IEEE Signal Processing Magazine*, vol. 19, no. 2, pp. 30–39, 2002.
- [13] Fredrik Gustafsson and Fredrik Gunnarsson, "Positioning using time-difference of arrival measurements," in *2003 IEEE International Conference on Acoustics, Speech, and Signal Processing, 2003. Proceedings.(ICASSP'03)*. IEEE, 2003, vol. 6, pp. VI–553.
- [14] Kenichi Kumatani, John McDonough, Jill Fain Lehman, and Bhiksha Raj, "Channel selection based on multichannel cross-correlation coefficients for distant speech recognition," in *2011 Joint Workshop on Hands-free Speech Communication and Microphone Arrays*. IEEE, 2011, pp. 1–6.
- [15] Martin Wolf and Climent Nadeu, "Channel selection measures for multi-microphone speech recognition," *Speech Communication*, vol. 57, pp. 170–180, 2014.
- [16] SM Dance and BM Shield, "The complete image-source method for the prediction of sound distribution in non-diffuse enclosed spaces," *Journal of Sound and Vibration*, vol. 201, no. 4, pp. 473–489, 1997.
- [17] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [18] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," *Advances in Neural Information Processing Systems*, vol. 33, pp. 12449–12460, 2020.
- [19] Tianyu Gao, Xingcheng Yao, and Danqi Chen, "SimCSE: Simple contrastive learning of sentence embeddings," in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, Online and Punta Cana, Dominican Republic, Nov. 2021, pp. 6894–6910, Association for Computational Linguistics.