

# FEDCPC: AN EFFECTIVE FEDERATED CONTRASTIVE LEARNING METHOD FOR PRIVACY PRESERVING EARLY-STAGE ALZHEIMER’S SPEECH DETECTION

Wenqing Wei<sup>1</sup>, Zhengdong Yang<sup>2,4</sup>, Yuan Gao<sup>2</sup>, Jiyi Li<sup>3</sup>, Chenhui Chu<sup>2</sup>, Shogo Okada<sup>1</sup>, Sheng Li<sup>4\*</sup>

<sup>1</sup>Japan Advanced Institute of Science and Technology, Nomi, Japan

<sup>2</sup>Kyoto University, Kyoto, Japan

<sup>3</sup>University of Yamanashi, Kofu, Japan

<sup>4</sup>National Institute of Information and Communications Technology, Kyoto, Japan

## ABSTRACT

The early-stage Alzheimer’s disease (AD) detection has been considered an important field of medical studies. Like traditional machine learning methods, speech-based automatic detection also suffers from data privacy risks because the data of specific patients are exclusive to each medical institution. A common practice is to use federated learning to protect the patients’ data privacy. However, its distributed learning process also causes performance reduction. To alleviate this problem while protecting user privacy, we propose a federated contrastive pre-training (FedCPC) performed before federated training for AD speech detection, which can learn a better representation from raw data and enables different clients to share data in the pre-training and training stages. Experimental results demonstrate that the proposed methods can achieve satisfactory performance while preserving data privacy.

**Index Terms** Alzheimer’s disease speech detection, federated learning, contrastive pre-training

## 1. INTRODUCTION

Alzheimer’s disease (AD) has been a global problem, and the number of people affected by this kind of cognitive impairment is growing. Unfortunately, this disease still cannot be perfectly cured. Therefore, early-stage AD detection methods are required. Current bio-medical diagnoses require a comprehensive examination by medical experts, which is costly and time-consuming. Compared with these biochemical methods, machine learning-based methods from easily captured spoken language signals are much more direct and efficient. Previous works have found speech changes in fluency [1], prosody [2], and rhythm [3, 4] in patients with AD. Researchers have been motivated to study AD detection using state-of-the-art technologies, such as speech recognition [5, 6, 7, 8], speaker recognition [9], natural language processing [10, 11, 12, 13], and multi-modeling [14].

Although such detection would be helpful, using patient data during model training might raise concerns about privacy issues [15, 16, 17]. The data of specific patients are exclusive for each medical institution, which cannot support traditional machine learning methods where all the local data are uploaded to one central server for learning a global model. Under such a setting, each medical institution does not have sufficient training data for the model; in addition, different medical institutions cannot benefit from a larger training data scale, leading to significant performance degradation. Existing methods focus on directly anonymizing speakers’ identities [18, 19, 20, 21, 22] and achieve very good results. However, these approaches are neither cheap nor time-efficient. Federated learning [23] is a machine learning technique that trains an algorithm across multiple decentralized edge devices or servers holding local data samples without the exchange of data to each other. It enables multiple clients to build a common and robust machine learning model without sharing data, thus preserving data privacy. Considering the advantages of the federated learning method, numerous works have recently used federated learning to protect user data privacy. Li et al. [24] proposed ADDETECTOR, a privacy-preserving smart healthcare system, to realize low-cost AD. Wang et al. [25] proposed a Blockchain-based Privacy-preserving Federated Learning scheme, which can enable the verifiability of the local models while protecting data privacy. However, federated learning is a kind of distributed learning, so the training data on each client will be smaller than in centralized learning, which leads to performance degradation, especially for small dataset tasks.

Several substantial works have recently demonstrated that self-supervised representations are highly successful in downstream speech and language processing tasks through feature-based speech representation extraction or fine-tuning as part of the downstream model [26]. Oord et al. [27] proposed contrastive predictive coding (CPC) that seeks to group samples that are alike while keeping samples that are different from one another apart from representation learning, which can accurately represent the data. Baevski et al. [28] used adversar-

\*Corresponding author. Wenqing Wei and Zhengdong Yang were NICT interns during this work.

ial training to train an unsupervised speech recognition model using the representations of the unlabeled speech audio data and the unlabeled phonemicized text data.

Inspired by these works, this paper proposes using federated contrastive learning to protect patient data privacy and enhance the model’s performance for each client by enabling data sharing while preserving privacy during the pre-training and training stages. Specifically, each medical institution could be regarded as a client. The clients locally train an independent contrastive predictive coding pre-training model and then upload the model to a central server. On the server side, the multiple models are then stacked up as a global model with federated averaging and hierarchical optimization, which cannot backtrack the parameters of individual models. After completing the Federated contrastive (FedCPC) pre-training process, we apply the FedCPC pre-training model as a feature extractor in each client to detect the AD speech with federated learning. Finally, this global model is sent to each client to benefit from more extensive data and guarantee strong anonymity and privacy.

## 2. METHODS

Figure 1 demonstrates the flowchart of the proposed method. Our method is divided into two steps. Step 1 is the proposed Federated CPC pre-training (FedCPC). The pre-training model for learning the representations of the speech data is CPC, and Federated learning is utilized for training the CPC model among the central server and the clients. After Federated CPC pre-training, Step 2 loads the pre-trained model and fine-tunes the classifier for the downstream Alzheimer’s disease detection task. The details of the CPC Model and the Federated Learning module are described as follows:

### 2.1. Contrastive Predictive Coding Pre-training

Self-supervised learning obtains supervisory signals from the data, which provides supervisory signals beneficial for downstream tasks. CPC [27] is a self-supervised learning approach, which involves predicting the future time step of data according to the context vector derived from past data. Its goal is to learn representations that allow long-term prediction of future time steps by maximizing mutual information between representations and predictions. Moreover, CPC captures high-level information from a signal (for example, global structures such as phonemes in speech [27]). This work adopts the CPC method for unsupervised representation pre-training of Alzheimer’s speech.

A convolutional encoder network produces a sequence of representation  $Z$  of a raw audio waveform. Subsequently, a recurrent context network summarizes the past information in the vector embedding sequence and produces corresponding contextual  $\mathbf{C} = \{\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_t\}$  representations for audio. The CPC pre-training model is trained to predict the future

latent representation  $\mathbf{z}_{t+k}$  ( $k$  is the predicted step) using the context-aware representation  $\mathbf{c}_t$  at the  $t$ -th timestep of speech. At each step  $t$ , we adopt a contrastive estimation-based InfoNCE [27] to maximize the mutual information lower bound between contextual representations  $\mathbf{c}_t$  and future latent representations  $\mathbf{z}_{t+k}$ . Here, given a set  $Z = \{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N\}$   $N$  random samples which contains one positive sample from  $p(\mathbf{z}_{t+k}|\mathbf{c}_t)$  and  $N - 1$  negative samples from “noise” distribution  $p(\mathbf{z}_{t+k})$  are drawn for optimizing the loss. And The formula is as follows:

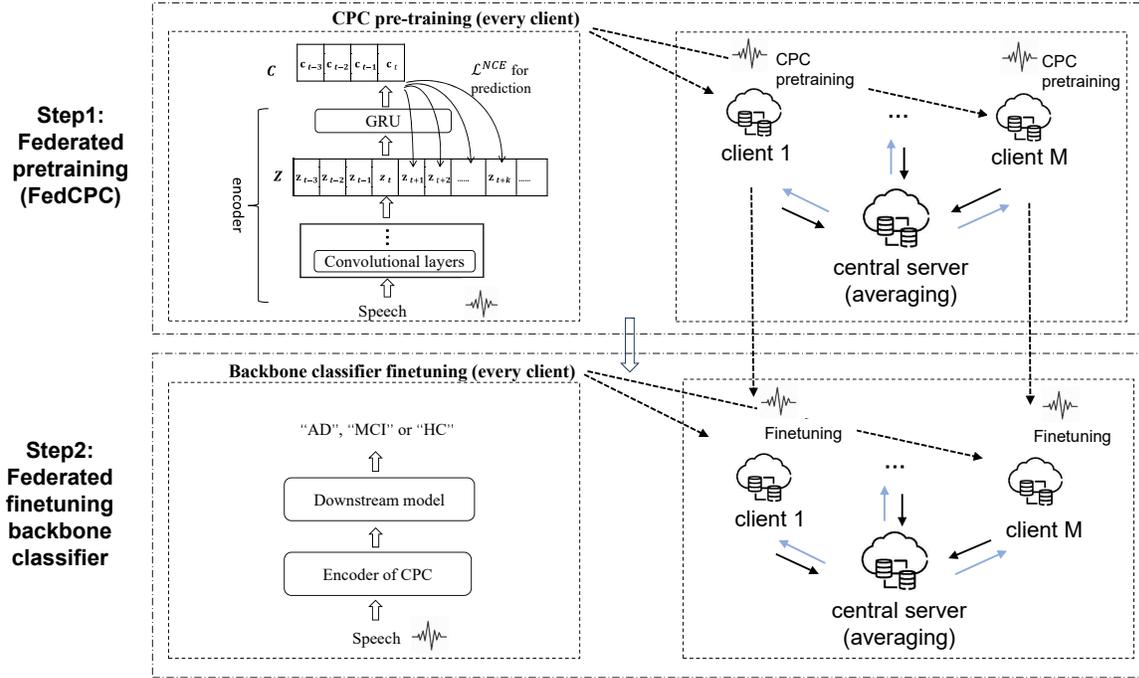
$$\mathcal{L}_{tk}^{NCE} = -\mathbb{E} \left[ \log \frac{\exp(\mathbf{c}_t^T \mathbf{W}_k \mathbf{z}_{t+k})}{\frac{1}{N} \sum_{\tilde{\mathbf{z}} \in Z} \exp(\mathbf{c}_t^T \mathbf{W}_k \tilde{\mathbf{z}})} \right] \quad (1)$$

### 2.2. Federated Learning

With the advent of the General Data Protection Regulation and increasing privacy concerns, sharing real-world AD speech data is facing significant challenges [15]. Speech data contain the speaker’s identifiable information represented as voiceprint used in many authentication systems [29, 30]. Exposing an individual’s voiceprint may cause security risks [16, 17] to voice authentication systems. Several methods have been proposed to anonymize speakers’ identities using state-of-the-art technologies, such as speech synthesis [18], voice conversion [19, 20], speaker embeddings [21], k-anonymity model [22]. However, these approaches are neither cost nor time-efficient and, therefore, do not meet the demands for big data on a global scale.

Federated learning proposes a solution to privacy-preserving large-scale machine learning by first training models locally on an individual client device and then aggregating the updates of the local models on the central server. There are several proposed algorithms for solving the Federated Optimization problem. One of the most promising algorithms is federated averaging (FedAvg) [23]. The algorithm provides privacy by design and personalizes models to individual users. It also claims to be more resource-efficient in terms of communication rounds.

This algorithm, firstly proposed in [23], relies on the stochastic gradient descent (SGD) optimization method since most of the most successful deep learning works were based on this. The available clients locally compute their average gradient on their local data at the parameters of the current model  $w_s$ , where  $s$  identifies the federated round. The central server aggregates these gradients and applies the update  $w_{s+1} \leftarrow w_s - \eta \sum_{m=1}^M \frac{n_m}{n} g_s^m$ , where  $g_s^m = \nabla F_m(w_s)$  is the gradient of the client  $m$  in  $s_{th}$  federated round,  $\eta$  is the learning rate,  $n_m$  is the number of samples at the client  $m$ ,  $n$  is the total number of samples (sum over all the available clients),  $M$  is the number of the clients. Equivalently, the update can be given by  $w_{s+1} \leftarrow \sum_{m=1}^M \frac{n_m}{n} w_{s+1}^m$ , where  $w_{s+1}^m \leftarrow w_s - \eta g_s^m, \forall m$ . Finally, every client takes a complete gradient descent step, while the server only takes the weighted average of the resulting models.



**Fig. 1.** The flowchart of the proposed method. The algorithm on the left side of the diagram will be applied to each client on the right side.

### 2.3. Overview of proposal model

As described above, federated learning is used in both steps of our proposed FedCPC-based model. For the FedPCPC pre-training model proposed in step 1,  $w_s$  contains all the parameters of the FedPCPC pre-training model in the  $s_{th}$  federated round. These parameters are shared through federated learning, which enables different clients to learn speech representation. For the downstream part of the FedCPC-based model of step 2,  $w_s$  contains the CPC pre-training model’s encoder part and the downstream model’s parameters in  $s_{th}$  federated round.

## 3. EXPERIMENTS

### 3.1. Data Description

The 2021 NCMMS AD Recognition Challenge provides the dataset. Since we joined the short speech track, we segmented all the long sentences from the data provided by the organizers into short speech clips of 6 seconds. As shown in Table 1, we selected a 7.16-hour speech set from 39 male speakers and 54 female speakers as the training set (Training) and a 0.67-hour

speech set from 15 male and 15 female speakers as the development set (Development). The test set is the official short speech track test set with a 1.92-hour speech set (Testing). All the sentences have labels in three kinds: Alzheimer’s disease (AD), mild cognitive impairment (MCI), and health common (HC).

### 3.2. System Implementation

We construct the following models both with conventional centralized training and federated learning, and their implementations are described as follows:

#### 3.2.1. Baseline

1. A 20-dimensional Mel frequency cepstral coefficients (MFCC) vector extracted with a 25ms window and 10ms frameshift is employed as the input feature of each frame. The CNN (convolutional neural network)-based system<sup>1</sup> has  $259 \times 20$  nodes as the input layer, three nodes as output, and five convolutional layers

<sup>1</sup>Available at <https://github.com/THUsatlab/AD2021>

**Table 1.** Data Descriptions

	#speakers (Male/Female)	#utterance (HC/MCI/AD)	#hour
Training	39 / 54	1712 / 1380 / 1208	7.16
Develop.	15 / 15	114 / 145 / 138	0.67
Testing	/	432 / 378 / 343	1.92

**Table 2.** Major Experimental Settings

<b>Training settings</b>	GPUs (3090)	1
	Batch-size	128
	Epochs	100
	Steps	34
	Optimizer	Adam
	Valid	Dev. Set
<b>Federated settings</b>	Federated round	50
	Local epochs	4
	Weighting Strategies	FedAvg
	#Clients	3

with one max-pooling layer following every layer. The convolutional filters of each layer are sequentially arranged as 32, 32, 32, 64, and 128. Finally, the last two layers are fully connected (FC) with 256 nodes before softmax output.

2. CNN is mainly for learning local features. To consider contextual dependencies from local features, we construct a CNN-LSTM-based model. Two LSTM layers are put on top of the CNN-based model, and this model also has two FC layers before softmax output. The network configuration of the CNN and FC layers is the same as the CNN-based system.
3. Moreover, to compare with other pre-trained models, we adopt AST<sup>2</sup> model [31]. It is transformer-based model for audio classification with a larger database with a simple architecture and superior performance. For this reason, we use the AST model to compare with our CPC-base model.

### 3.2.2. Self-Supervised Model

**Pre-training model setup:** The implementation of CPC pre-training model is similar to [27]. Each training iteration randomly extracts a segment of about 20,480 frames from every utterance as the encoder’s input. The encoder comprises five 1-dimensional CNN layers and a single-layer gated recurrent unit (GRU). In detail, each of the five layers has the same down-sampling rate of 1/160 to get the same frame rate and the same settings of the filter size, strides, and paddings ([10, 8, 4, 4, 4], [5, 4, 2, 2, 2] and [3, 2, 1, 1, 1]). All five layers

have 512 hidden units. Moreover, the GRU layer is employed as the sequence model with 256 hidden units. Every frame of GRU output is used to predict the context  $c$  (12 future frames). Adam optimizer trains the model with a learning rate of  $2e-4$  and a minibatch size 8.

**Backbone classifier network:** As AST is a pre-trained model used to compare our proposed method, we only constructed CNN and CNN-LSTM models as downstream models in this work. The training parameters are the same as the baseline and the scheme above.

We used the open-sourced federated learning framework Flower [32] and PyTorch version-1.9.1<sup>3</sup> for all our experiments. Table 2 lists the training and test settings. Moreover, speakers of each client’s data are unique to ensure the independence of speakers from different clients in federated learning.

## 4. REUSLTS AND DISCUSSIONS

### 4.1. Main results

Using the centralized learning scenario, we first investigate whether the CPC pre-training model can improve performance. Table 3 shows the experimental results of the test set of centralized learning. Compared with the CNN model, the CNN-LSTM model shows better performances in both machine learning paradigms, especially outperforming the CNN system on precision, recall, and F1-score by 3.8%, 4.2%, and 5.0% in centralized learning. Meanwhile, using CPC pre-training models leads to better performance than MFCC features. As shown in Table 3, the macro F1-score achieved 8.5%, 3.6% improvements in centralized learning by comparing CNN and CNN-LSTM models, and CPC-CNN almost got the same F1-score with CPC-LSTM-CNN. Moreover, the paper [8] utilizes the Wav2vec speech recognition pre-training model to predict Alzheimer’s Disease on the same databases. In contrast, our proposed method leverages the CPC pre-training model without federated learning, specifically the CPC-CNN-LSTM model, and achieves a superior result (F1: 78.8%) compared to the best result obtained by Wav2vec2.0.3-2 (F1: 77.2%) in [8] on short audio tracks. These findings provide evidence that the CPC pre-training model captures structural information embedded in raw audio signals, enhancing AD detection performance. Figure 2 illustrates the visualization of representations learned by the pre-training model, indicating specific properties that can be utilized for clustering in this task.

We then verify the performance of our proposed FedCPC-based models. Table 4 demonstrates the results of federated learning. The CPC-Client-based model used the pre-training models trained by different clients using data unique to each client. The FedCPC-based model used the pre-training models trained with federated learning. As shown in Table 4, the

<sup>2</sup>Available at <https://github.com/YuanGongND/ast>

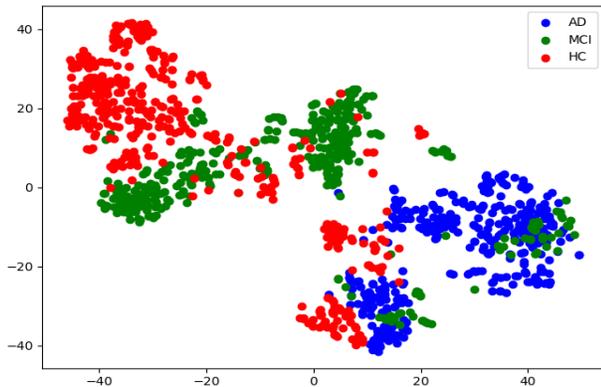
<sup>3</sup><https://pytorch.org/>

**Table 3.** Evaluation with centralized learning.

	Precision(%)	Recall(%)	F1(%)
CNN	71.3	71.1	70.2
CNN-LSTM	75.1	75.3	75.2
CPC-CNN	80.3	78.4	78.7
CPC-CNN-LSTM	<b>84.7</b>	<b>78.4</b>	<b>78.8</b>
AST	75.2	74.3	75.5
Wav2vec2.0_3-2 [8]	77.9	77.2	77.2

**Table 4.** Evaluation on the testing set with federated learning.

Method	Downstream Model	Precision(%)	Recall(%)	F1(%)
Fed-AST	-	73.7	72.5	73.5
Fed	CNN	72.8	70.2	68.2
CPC-Client	CNN	77.0	71.8	72.3
FedCPC (Our)	CNN	<b>79.3</b>	<b>74.3</b>	<b>74.8</b>
Fed	CNN-LSTM	72.0	72.1	71.4
CPC-Client	CNN-LSTM	76.4	72.9	73.5
FedCPC (Our)	CNN-LSTM	<b>79.5</b>	<b>74.7</b>	<b>75.3</b>

**Fig. 2.** t-SNE visualization of audio with CPC pre-training.

FedCPC-based model outperforms the CPC-Client model. Our proposed FedCPC-CNN-LSTM model gets the best result in federated learning. Compared with the CNN-LSTM model, the FedCPC-CNN-LSTM model significantly improves from 71.4% to 75.3 % in the F1-score. For comparison, we also exploit a pre-training AST model on this task and notice that FedCPC-CNN-LSTM outperforms it. Meanwhile, we check the results for different categories by using confusion matrices to explore the reasons for the improvement. The confusion matrices for Fed-CNN-LSTM and FedCPC-CNN-LSTM with features computed on the test set, as depicted in Figure 3. We found that Fed-CNN-LSTM tends to classify AD into the MCI and HC category incorrectly, while FedCPC-CNN-LSTM tends to classify AD into the MCI. It is consistent with the distribution of audio in Figure 2. Moreover, FedCPC-CNN-LSTM performs better on AD and MCI.

True Labels	Predicted Labels		
	AD	MCI	HC
AD	173	82	88
MCI	69	292	17
HC	5	62	365

a) Fed-CNN-LSTM

True Labels	Predicted Labels		
	AD	MCI	HC
AD	231	112	0
MCI	45	325	8
HC	15	112	305

b) FedCPC-CNN-LSTM with Federated learning

**Fig. 3.** Confusion matrix of federated learning (a) Fed-CNN-LSTM, (b) FedCPC-CNN-LSTM.

## 4.2. Further Discussions

With the results in Tables 3 and 4, we found that the results of the federated-learning-based models have decreased compared with the traditional centralized learning results. The recall and F1-score of the CNN-LSTM system decreased by 3.1%, 3.2%, and 3.8% after adopting federated learning. Moreover, the FedCPC-CNN-LSTM dropped by 1.8% and 3.4% in the recall and F1-score after adopting federated learning, respectively. This result is expected regarding each client side’s limited data and partial observations. This shows the inherent limitations of federated learning, which are worth investigating. The AST cannot outperform the pre-trained small models in the experiments above. Considering its high cost, it is not applicable for memory-restricted clients.

## 5. CONCLUSION

This paper uses federated learning to train models to detect AD in speech and protect privacy in the user’s voice. Compared with centralized learning, federated learning is a distributed training mode, leading to performance degradation, especially in small databases. To address this issue, we proposed the FedCPC-based pre-training method, which enables data sharing while preserving privacy during the pre-training and training stage. The experimental evaluation revealed that our proposed approach effectively preserves privacy while maintaining competitive performance compared to non-privacy-preserving methods.

## 6. ACKNOWLEDGEMENT

This work was supported by JSPS KAKENHI Grant Numbers JP23K11227, JP23H03454, and NICT tenure-track funding.

## 7. REFERENCES

- [1] Greta Szatloczki, Ildiko Hoffmann, Veronika Vincze, Janos Kalman, and Magdolna Pakaski, “Speaking in alzheimer’s disease, is that an early sign? importance of changes in language abilities in alzheimer’s disease,” *Frontiers in aging neuroscience*, vol. 7, pp. 195, 2015.

- [2] Kaye Horley, Amanda Reid, and Denis Burnham, "Emotional prosody perception and production in dementia of the alzheimer's type," 2010.
- [3] Anja Lowit, Bettina Brendel, Corinne Dobinson, and Peter Howell, "An investigation into the influences of age, pathology and cognition on speech production," *Journal of medical speech-language pathology*, vol. 14, pp. 253, 2006.
- [4] Brian Roark, Margaret Mitchell, John-Paul Hosom, Kristy Hollingshead, and Jeffrey Kaye, "Spoken language derived measures for detecting mild cognitive impairment," *IEEE transactions on audio, speech, and language processing*, vol. 19, no. 7, pp. 2081–2090, 2011.
- [5] László Tóth, Gábor Gosztolya, Veronika Vincze, Ildikó Hoffmann, Gréta Szatlóczi, Edit Biró, Fruzsina Zsura, Magdolna Pákáski, and János Kálmán, "Automatic detection of mild cognitive impairment from spontaneous speech using asr," in *INTERSPEECH*, 2015.
- [6] Yilin Pan, Bahman Mirheidari, Markus Reuber, Annalena Venneri, Daniel J. Blackburn, and Heidi Christensen, "Improving detection of alzheimer's disease using automatic speech recognition to identify high-quality segments for more robust feature extraction," in *INTER-SPEECH*, 2020.
- [7] Luke Zhou, Kathleen C. Fraser, and Frank Rudzicz, "Speech recognition in alzheimer's disease and in its assessment," in *INTERSPEECH*, 2016.
- [8] Ying Qin, Wei Liu, Zhiyuan Peng, Si-Ioi Ng, Jingyu Li, Haibo Hu, and Tan Lee, "Exploiting pre-trained asr models for alzheimer's disease recognition through spontaneous speech," *arXiv preprint arXiv:2110.01493*, 2021.
- [9] Raghavendra Reddy Pappagari, Jaejin Cho, Laureano Moro-Velázquez, and Najim Dehak, "Using state of the art speaker recognition and natural language processing technologies to detect alzheimer's disease and assess its severity," in *INTERSPEECH*, 2020.
- [10] Thomas Searle, Zina M. Ibrahim, and Richard J. B. Dobson, "Comparing natural language processing techniques for alzheimer's dementia prediction in spontaneous speech," in *INTERSPEECH*, 2020.
- [11] Sebastian Wankerl, Elmar Nöth, and Stefan Evert, "An n-gram based approach to the automatic diagnosis of alzheimer's disease from spoken language," in *INTERSPEECH*, 2017.
- [12] Aparna Balagopalan, Benjamin Eyre, Frank Rudzicz, and Jekaterina Novikova, "To bert or not to bert: Comparing speech and language-based approaches for alzheimer's disease detection," in *INTER-SPEECH*, 2020.
- [13] Jiahong Yuan, Yuchen Bian, Xingyu Cai, Jiaji Huang, Zheng Ye, and Kenneth Church, "Disfluencies and fine-tuning pre-trained language models for detection of alzheimer's disease," in *INTERSPEECH*, 2020.
- [14] Morteza Rohanian, J. Hough, and Matthew Purver, "Multi-modal fusion with gating using audio, lexical and disfluency features for alzheimer's dementia recognition from spontaneous speech," *ArXiv*, vol. abs/2106.09668, 2020.
- [15] A. Nautsch and et al., "The GDPR & speech data: Reflections of legal and technology communities, first steps towards a common understanding," 2019.
- [16] Z. Wu and et al., "Spoofing and countermeasures for speaker verification: A survey," *Speech Communication*, vol. 66, pp. 130–153, 2015.
- [17] S. Suwajanakorn and et al., "Synthesizing obama: learning lip sync from audio," *ACM Transactions on Graphics*, vol. 36, no. 4, pp. 95, 2017.
- [18] T. Justin and et al., "Speaker de-identification using diphone recognition and speech synthesis," in *11th IEEE International Conference and Workshops on Automatic Face and Gesture Recognition*. IEEE, 2015, vol. 4, pp. 1–7.
- [19] J. Qian and et al., "Hidebehind: Enjoy voice input with voiceprint unclonability and anonymity," in *Proceedings of the 16th ACM Conference on Embedded Networked Sensor Systems*. ACM, 2018, pp. 82–94.
- [20] Brij Mohan Lal Srivastava, Nathalie Vauquier, Md Sahidullah, Aurélien Bellet, Marc Tommasi, and Emmanuel Vincent, "Evaluating voice conversion-based privacy protection against informed attackers," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 2802–2806.
- [21] F. Fang and et al., "Speaker anonymization using X-vector and neural waveform models," in *Proc. 10th ISCA Speech Synthesis Workshop*, 2019, pp. 155–160.
- [22] L. Sweeney, "K-anonymity: A model for protecting privacy," *Int. J. Uncertain. Fuzziness Knowl.-Based Syst.*, vol. 10, no. 5, pp. 557–570, 2002.
- [23] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Artificial intelligence and statistics*. PMLR, 2017, pp. 1273–1282.
- [24] Jiachun Li, Yan Meng, Lichuan Ma, Suguo Du, Haojin Zhu, Qingqi Pei, and Xuemin Shen, "A federated learning based privacy-preserving smart healthcare system," *IEEE Transactions on Industrial Informatics*, vol. 18, no. 3, 2021.
- [25] Yangpeng Wang, Ling Xiong, Xianhua Niu, Yunxiang Wang, and Dexin Liang, "A federated learning based privacy-preserving data sharing scheme for internet of vehicles," in *Frontiers in Cyber Security: 5th International Conference, FCS 2022, Kumasi, Ghana, December 13–15, 2022, Proceedings*. Springer, 2022, pp. 18–33.
- [26] Ashish Jaiswal and et al., "A survey on contrastive self-supervised learning," *Technologies*, vol. 9, no. 1, pp. 2, 2020.
- [27] Aaron van den Oord, Yazhe Li, and Oriol Vinyals, "Representation learning with contrastive predictive coding," *arXiv preprint arXiv:1807.03748*, 2018.
- [28] Alexei Baevski, Wei-Ning Hsu, Alexis Conneau, and Michael Auli, "Unsupervised speech recognition," *Advances in Neural Information Processing Systems*, vol. 34, pp. 27826–27839, 2021.
- [29] A. Boles and P. Rad, "Voice biometrics: Deep learning-based voiceprint authentication system," in *12th System of Systems Engineering Conference*, 2017, pp. 1–6.
- [30] Tencent Inc., "The new wechat password," <https://blog.wechat.com/tag/voiceprint/>, 2015.
- [31] Yuan Gong, Yu-An Chung, and James Glass, "AST: Audio Spectrogram Transformer," in *Proc. Interspeech 2021*, 2021, pp. 571–575.
- [32] Daniel J Beutel, Taner Topal, Akhil Mathur, Xinchu Qiu, Javier Fernandez-Marques, Yan Gao, Lorenzo Sani, Kwing Hei Li, Titouan Parcollet, Pedro Porto Buarque de Gusmão, et al., "Flower: A friendly federated learning framework," 2022.