# Clever Hans Effect Found in Automatic Detection of Alzheimer's Disease through Speech

*Yin-Long Liu[1], Rui Feng[1], Jiahong Yuan[1,2*], Zhen-Hua Ling[1,2]*

[1]National Engineering Research Center of Speech and Language Information Processing,
University of Science and Technology of China, Hefei, P. R. China
[2]Interdisciplinary Research Center for Linguistic Sciences,
University of Science and Technology of China, Hefei, P. R. China

{lyl2001, fengruimse}@mail.ustc.edu.cn, {jiahongyuan, zhling}@ustc.edu.cn

## Abstract

We uncover an underlying bias present in the audio recordings produced from the picture description task of the Pitt corpus, the largest publicly accessible database for Alzheimer's Disease (AD) detection research. Even by solely utilizing the silent segments of these audio recordings, we achieve nearly 100% accuracy in AD detection. However, employing the same methods to other datasets and preprocessed Pitt recordings results in typical levels (approximately 80%) of AD detection accuracy. These results demonstrate a Clever Hans effect in AD detection on the Pitt corpus. Our findings emphasize the crucial importance of maintaining vigilance regarding inherent biases in datasets utilized for training deep learning models, and highlight the necessity for a better understanding of the models' performance.

**Index Terms**: Alzheimer's disease detection, spurious features, bias, Clever Hans effect

## 1. Introduction

Alzheimer's Disease (AD), the most common cause of dementia, is a neurodegenerative disease that worsens over time and causes irreversible damage to the brain, manifested by a persistent deterioration of an individual's cognitive and functional abilities, including language, memory, attention, and executive function [1].

In recent years, researchers have achieved promising results in utilizing deep-learning models and an end-to-end approach for the automatic detection of AD through speech. However, the robustness of these models have not been thoroughly tested, primarily due to the scarcity of large and diverse datasets.

In this paper, we reveal an inherent bias present in the audio recordings produced from the picture description task of the Pitt corpus from DementiaBank, the largest publicly accessible database for AD detection research. Remarkably, even by exclusively utilizing the silent segments of these audio recordings, we achieve nearly 100% accuracy in AD detection. As far as we are aware, this bias has not been reported in the literature.

We present this finding to draw researchers' attention to the impact of bias in the dataset, and advocate for more effort in studying the robustness and explainability of deep-learning models in automatic AD detection.

## 2. Related work

### 2.1. The Pitt corpus

The Pitt corpus [2] is a widely used subset of DementiaBank. It was collected over a longitudinal period, encompassing 104
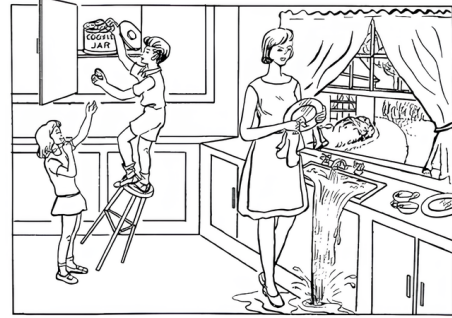
Figure 1: *The picture of "Cookie Theft", adopted from Boston Diagnostic Aphasia Examination.*

elderly controls, 208 individuals with probable or possible AD, and 85 participants with unknown diagnoses. Responses to four language tasks were recorded, including one task of describing the content of the Cookie Theft picture for all participants, which was originally designed for the Boston Diagnostic Aphasia Examination [3] (as shown in Figure 1), and three tasks of verbal fluency, sentence construction and story recall for AD participants only. For picture description task, there are 306 AD speech samples, 243 Healthy Controls (HC) speech samples, and due to the interference of the recording environment, these speech samples contain noise.

### 2.2. Automatic detection of AD through speech

Many researchers have utilized the Pitt corpus for AD detection studies. As of now, according to Google Scholar, this dataset has been cited in 543 papers, and the citation count is increasing year by year.

Table 1 presents some results of previous studies using features extracted from the original, non-denoised recordings of the Pitt corpus. For example, Han et al. [4] utilized a miniature version of the Xception network [5] to develop a deep learning model classifying AD and HC based on original speech selected from the Pitt corpus using log-mel spectrogram features, achieving an accuracy of 94.2%.

The Pitt corpus has been used for AD detection chal-

Table 1: *Some results of previous studies using features extracted from the original, non-denoised recordings of the Pitt corpus.*

| References | Speech from Pitt | Results(%) |
|---|---|---|
| Han et al. [4] | 217 AD, 242 HC | 94.2 Accuracy |
| Ammar et al. [6] | 43 AD, 43 HC | 91 F-measure |
| Zargarbashi et al. [7] | 255 AD, 233 HC | 83.6 Accuracy |
| Fraser et al. [8] | 240 AD, 233 HC | 81.92 Accuracy |

Table 2: *The accuracy of previous studies using only audio recordings from challenge datasets for AD detection.*

| References | Dataset | Accuracy(%) |
|---|---|---|
| Mei et al. [11] | ADReSS | 79.2 |
| Koo et al. [12] | ADReSS | 72.9 |
| Gauder et al. [13] | ADReSSo | 78.9 |
| Chen et al. [14] | ADReSSo | 77.1 |

lenges at international conferences, including Interspeech 2020 (ADReSS) [9] and Interspeech 2021 (ADReSSo) [10]. In these challenges, the original recordings of the Pitt corpus underwent denoising and normalization process to create training and test data. Table 2 presents the accuracy of previous studies using only audio recordings from challenge datasets for AD detection. As we can see, compared to the original, non-denoised data, the AD detection accuracy on the preprocessed data from these challenges was notably lower.

The performance gap between using the original and preprocessed Pitt corpus may be attributed to two potential reasons. One is the improvement of the model's capability, and the other is the impact of spurious features present in the original speech data, i.e., the Clever Hans effect.

## 2.3. The clever hans effect

During the optimization process, models may exploit spurious correlations in the training data, resulting in seemingly high-performance metrics, a phenomenon known as the Clever Hans effect [15]. The Clever Hans refers to a horse that was believed to perform arithmetic and other intellectual activities. Subsequent investigations revealed that the horse did not actually execute these intellectual tasks. Instead, it responded to involuntary cues in the body language of its human trainer, of which the human trainer was entirely unaware.

The Clever Hans effect is frequently observed in the context of supervised classifiers. Several machine learning problems have illustrated the effect. Arjovsky et al. [16] trained a convolutional neural network designed to classify camels and cows. After experimental analysis, it was discovered that the neural network had successfully minimized its training error through a simple cheat: categorizing green landscapes as cows and beige landscapes as camels. Borah et al. [17] mentioned Clever Hans behavior in high-performance neural translationese classifiers, where BERT-based classifiers capitalize on spurious correlations, in particular topic information, between data and target classification labels, rather than genuine translationese signals. Chettri et al. [18] proposed that any visible pattern difference, such as the distribution of silence, between bonafide and spoof classes can introduce biases in voice spoofing detection, consequently influencing model decisions. Similar effect has also been demonstrated in medical contexts. Wallis et al. [19] exposed an underlying bias in a commonly used publicly available brain tumour MRI dataset, and proposed that this is due to implicit radiologist input in the selection of the 2D slices. In the KDD CUP breast cancer identification challenge, Perlich etal. [20] found that the patient IDs (which had not been removed from the data) were highly correlated with the malignancy of the patients' tumours. Several recent studies [21, 22] have revealed biases in datasets designed for COVID-19 identification from X-ray images, which stem from the inclusion of positive and negative images obtained from distinct sources. The biases described in the above-mentioned studies can lead machine learning models to take a "shortcut" and address a significantly easier task.

# 3. Unveiling the Clever Hans effect

## 3.1. Data

In addition to employing speech recordings from the Pitt corpus, we also utilized Mandarin speech samples from iFLYTEK, as well as speech samples from ADReSS and ADReSSo, for comparative analysis.

### 3.1.1. Data from the original Pitt corpus

We selected 255 speech samples from 168 probable or possible participants and 242 speech samples from 99 HC participants. In addition to the Pitt corpus original (*Pco*, original speech) dataset, based on the timestamp information in manual transcripts of the corresponding speech recordings, we also derived these three datasets: Pitt corpus subject (*Pcsu*, containing only subject speech), Pitt corpus silence (*Pcsi*, containing only silent speech), and Pitt corpus interviewer (*Pci*, containing only interviewer speech). Given that a certain participant may have multiple corresponding speech samples, we randomly selected one of them to ensure that speech samples from a specific participant do not simultaneously appear in both the training and test datasets.

### 3.1.2. Data from a Mandarin AD dataset

The Mandarin data utilized in this paper is sourced from iFLYTEK, where subjects were recruited from the Department of Neurology and the Department of Memory Clinic of Shanghai Tongji Hospital [23, 24, 25] and were instructed to undertake the same picture description task mentioned earlier. We selected 120 AD speech samples and 173 HC speech samples. Similar to Section 3.1.1, we obtained four datasets: Mandarin original (*Mo*), Mandarin subject (*Msu*), Mandarin silence (*Msi*), and Mandarin interviewer (*Mi*). Each speech sample corresponds to a unique participant.

### 3.1.3. Data from the ADReSS and ADReSSo challenges

ADReSS and ADReSSo challenges were hosted by Interspeech 2020 [9] and Interspeech 2021 [10] conferences respectively. Both challenge datasets are in English and underwent acoustic enhancement through noise removal and audio volume normalization. The ADReSS dataset contains 78 AD speech samples and 78 HC speech samples respectively. The ADReSSo dataset contains 122 AD speech samples and 115 HC speech samples. In addition to the ADReSS original (*Ao*) and ADReSSo original (*Aoo*) datasets, we also derived ADReSS silence (*As*) and ADReSSo silence (*Aos*) datasets using the pyannote[1] voice activity detection (VAD) tools [26], since the timestamp information provided by these two challenge datasets does not include silent intervals.

## 3.2. Fine-tuning wav2vec 2.0 for AD detecion

Wav2vec 2.0 is a framework for self-supervised learning of speech representations using contrastive loss [27]. In previous work [28], we had demonstrated the effectiveness of fine-tuning wav2vec 2.0 for AD detection. In this paper, we fine-tuned the wav2vec 2.0 models "facebook/wav2vec2-large-xlsr-53" and "wbbbbb/wav2vec2-large-chinese-zh-cn" with a sequence classification head on top (a linear layer with the sigmoid activation function over the average pooled output) on English and Mandarin speech data respectively. The models are available in the HuggingFace's Transformers library[23]. We used 5-fold cross-

---

[1]https://github.com/pyannote/pyannote-audio
[2]https://huggingface.co/facebook/wav2vec2-large-xlsr-53
[3]https://huggingface.co/wbbbbb/wav2vec2-large-chinese-zh-cn

Table 3: *The AD detection accuracy of using different datasets to fine-tune wav2vec 2.0 models respectively.*

| Dataset | Different subdataset | Number of training/test samples | Accuacy(%) |
|---|---|---|---|
| Pitt corpus | *Pco* | 204/51 | 97.2 |
| | *Pcsu* | 124/31 | 90.3 |
| | *Pcsi* | 208/52 | 98.9 |
| Mandarin | *Mo* | 232/58 | 81 |
| | *Msu* | 136/34 | 73.5 |
| | *Msi* | 52/13 | 57.3 |
| ADReSS | *Ao* | 112/28 | 80.7 |
| | *As* | 120/30 | 56.7 |
| ADReSSo | *Aoo* | 176/44 | 77.7 |
| | *Aos* | 188/47 | 61.3 |

validation to evaluate the models.

For fine-tuning the wav2vec 2.0 models, we set the batch size to 1, the gradient accumulation steps to 4, the number of training epochs to 15, the learning rate to $3 \times 10^{-5}$, the warmup ratio to 0.1, and the loss function was cross-entropy. we employed the *Transformers.Trainer* as the optimizer. We converted the audio file format from mp3 to wav and converted the audio from stereo to mono, along with downsampling the audio data from 44.1kHz to 16kHz.

### 3.3. Results

The following three experiments progressively introduce how we unveil the Clever Hans effect.

Initially, we treated the Pitt corpus as a normal dataset for AD detection research. We used *Pco*, *Pcsu*, *Mo*, *Msu*, *Ao*, *Aoo* to fine-tune the wav2vec 2.0 models respectively. Only speech recordings with a duration longer than 35 seconds were retained in the datasets, and only the first 35 seconds were used for fine-tuning. The results are shown in the corresponding rows of Table 3. It can be seen that the accuracy of Pitt corpus is much higher than that of the other ones. Specifically, the classification accuracy of *Pco* is much higher than that of the two challenge datasets *Ao* (80.7%), *Aoo* (77.7%) and the Mandarin dataset *Mo* (81%) and it is close to 100% (97.2%), a result that is worth pondering, since it should be comparable to the performance on *Ao* and *Aoo*. Considering the above results, we initially suspect that the speech recordings in the Pitt corpus are interfered by some factors.

Next, for proving that the Pitt corpus indeed has problems, we conducted the second experiment. We fine-tuned the wav2vec 2.0 models using the first 85 seconds of each training sample in the two datasets, *Pcsu* and *Msu*, respectively, and conducted the test on the first 85 seconds of each speech sample from *Pci* and *Mi* datasets, respectively. The label of each speech sample in *Pci* and *Mi* is the same as the subject interviewed by the corresponding interviewer. The results are shown in Table 4. The test performance on the *Pci* dataset is an astonishing 83.1%, a figure that seems unbelievable given that the model, trained exclusively on subjects' speech, theoretically lacks the ability to identify a subject's AD based solely on the interviewer's speech. The test performance on the *Mi* dataset is relatively low, only 64.1%, which is normal. Based on the above results, we can confirm that the Pitt corpus is definitely influenced by certain factors.

Then, in order to further analyze which specific factors in-

Table 4: *The AD detection accuracy of the wav2vec 2.0 models fine-tuned with only subject speech on only subject speech or only interviewer speech respectively.*

| Training Set | Test Set | Accuacy(%) |
|---|---|---|
| *Pcsu* | *Pcsu* | 98.1 |
| | *Pci* | 83.1 |
| *Msu* | *Msu* | 84.4 |
| | *Mi* | 64.1 |

terfere with the Pitt corpus , we attempted the third experiment. We used the first 35 seconds of each speech sample from *Pcsi*, *As*, *Aos*, and *Msi* to fine-tune the wav2vec 2.0 models, respectively. We didn't use the duration information of the silent segments, instead, all silent segments of each original speech sample were concatenated into one piece as input for fine-tuning. The label for the silence piece is the same as the corresponding subject. The results are shown in the corresponding rows of Table 3. It can be seen that the performance on *Pcsi* can reach 98.9%, which is astonishing. On the contrary, the accuracy on *As* (56.7%) and *Aos* (61.3%) is much lower, and the accuracy on *Msi* (57.3%) is almost random guessing, as what we can expect. These results suggest that the audio recordings in Pitt corpus are interfered by environmental factors such as background noise. The models learned to capture these spurious features and correlations, leading to their high performance.

## 4. Validating the Clever Hans effect

### 4.1. Classification based on hand-crafted and wav2vec 2.0 features

To validate the bias of the recording environment in speech samples from the Pitt corpus, we employed the openSMILE toolkit [29] to extract the ComParE 2016 features [30] from speech recordings containing only silence, serving as our low-level acoustic features. ComParE 2016 is the largest feature set (6373 dimensions) in the toolkit and has been used for AD detection [10, 28, 14]. We utilized XGBoost, GBDT, AdaBoost classifiers, and their majority voting, as provided by the scikit-learn package, to conduct 5-fold cross-validation on the aforementioned features for both the English silent speech dataset from the Pitt corpus and Mandarin silent speech dataset. Likewise, these classifiers have also been successfully applied to AD detection [28, 31]. In addition, features from the last hidden layer of the wav2vec 2.0 model fine-tuned with the English silent speech dataset (*Pcsi*) were also utilized for building classifiers. Both the original 1024-dimensional features and dimensionally reduced ones were explored. The reduction was achieved to 10 and 5 dimensions using Principal Component Analysis (PCA) from the scikit-learn package.

We used the aforementioned methods to study *Pcsi* and *Msi* to validate our hypothesis, by building classifiers on the ComParE 2016 feature set for these two datasets. Additionally, we explored the 1024-dimensional features generated by the fine-tuned wav2vec 2.0 model with the *Pcsi* dataset, as well as the

Table 5: *AD detection accuracy (%) of each machine learning classifier on the ComParE 2016 feature set of the Pcsi and Msi datasets.*

| Dataset | XGBoost | GBDT | AdaBoost | Voting |
|---|---|---|---|---|
| *Pcsi* | 83.8 | 80.1 | 87.6 | 80.8 |
| *Msi* | 49.3 | 50.7 | 55.5 | 55.2 |

Table 6: *AD detection accuracy (%) of each machine learning classifier on the features generated by the fine-tuned wav2vec 2.0 model with Pcsi dataset and the fine-tuned wav2vec 2.0 model's dimensionally reduced features.*

| Feature set | XGBoost | GBDT | AdaBoost | Voting |
|---|---|---|---|---|
| wav2vec 2.0 (1024) | 97.4 | 98.9 | 98.5 | 98.9 |
| PCA(10) | 99.2 | 97.7 | 97.7 | 97.7 |
| PCA(5) | 99.2 | 98.1 | 98.1 | 98.1 |

dimensionally reduced features (10 dimensions and 5 dimensions) of the fine-tuned wav2vec 2.0 model. The results are shown in Table 5 and Table 6. It can be observed that AdaBoost classifier can still achieve a high AD detection accuracy of 87.6% based solely on *Pcsi*'s low-level acoustic features ComParE 2016, while the performance of each classifier on *Msi*'s ComParE 2016 is only about 50%. Moreover, whether it is on the original features (1024 dimensions) of the fine-tuned wav2vec 2.0 model with *Pcsi* or on the features after dimensionality reduction, we can achieve an unbelievable nearly 100% accuracy on the *Pcsi* dataset, which contain only silences of the Pitt corpus. These results confirm an underlying bias in the speech recordings produced by the picture description task in the Pitt corpus, which is caused by interference from environmental factors such as background noise. Studies using acoustic features extracted from the original speech recordings of this dataset for AD detection will be affected by this bias.

Figure 2 depicts the spectrograms of two randomly selected silent segments (AD and HC) from *Pcsi*. We can easily distinguish between AD and HC based on the two segments. This observation is consistent with the point made above.

### 4.2. Results from using preprocessed speech recordings

In order to mitigate the impact of the bias, we preprocessed the speech recordings before employing machine learning methods on the datasets. We utilized noisereduce package[4] to reduce stationary noise [32]. The package relies on a method called "spectral gating" which is a form of Noise Gate. It works by computing a spectrogram of a signal (and optionally a noise signal) and estimating a noise threshold (or gate) for each frequency band of that signal/noise. That threshold is used to compute a mask, which gates noise below the frequency-varying threshold. For stationary noise reduction, it should keep the estimated noise threshold at the same level across the whole signal. We also utilized the "AudioSegment" and "effects" methods of pydub
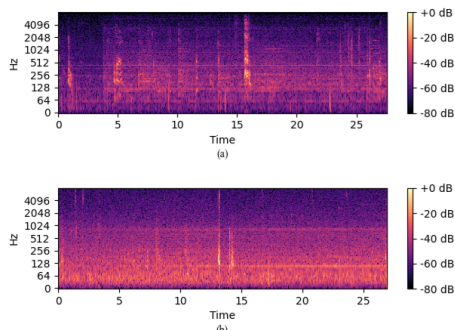


Figure 2: *Spectrograms of (a) an AD speech sample and (b) an HC speech sample in Pcsi.*

---

[4]https://github.com/timsainb/noisereduce

Table 7: *The AD detection accuracy of using preprocessed Pitt and Mandarin datasets to fine-tune wav2vec 2.0 models. (The numbers in the parentheses correspond to results obtained using original Pitt and Mandarin datasets listed in Table 3.)*

| Dataset | Different subdataset | Number of training/test samples | Accuacy(%) |
|---|---|---|---|
| Pitt corpus | *Pco* | 204/51 | 82 (97.2) |
| | *Pcsu* | 124/31 | 77.4 (90.3) |
| | *Pcsi* | 208/52 | 63.1 (98.9) |
| Mandarin | *Mo* | 232/58 | 80.7 (81) |
| | *Msu* | 136/34 | 74.8 (73.5) |
| | *Msi* | 52/13 | 58.8 (57.3) |

package[5] for standard amplitude normalization, which scaled the whole audio to the max amplitude.

We utilized the methods described above to preprocess the *Pco*, *Pcsu*, *Pcsi*, *Mo*, *Msu*, and *Msi*, respectively. Subsequently, under the same experimental parameters configuration as the unpreprocessed datasets, we employed the preprocessed new speech datasets to fine-tune the wav2vec 2.0 models respectively. The results are presented in Table 7. Comparing these outcomes with the corresponding results on original datasets, it is evident that the performance of the preprocessed *Pcsi* dataset has significantly decreased (from 98.9% to 63.1%). This reduced performance is now comparable to that of *Aos* (61.3%). Moreover, the performance on the preprocessed *Msi* remains at the level of random guessing (58.8%). Additionally, the performance on preprocessed *Pco* and *Pcsu* is reduced to what we consider a normal level. Specifically, the accuracy on *Pco* decreases from 97.2% to 82%, which is essentially similar to the performance on *Ao* (80.7%). The results further confirm that speech recordings in the original Pitt corpus are affected by environmental interference, such as background noise. The performance of preprocessed *Mo* and *Msu* is basically equivalent to or slightly improved compared to that of the original data, indicating the preprocessing methods do not compromise the valuable information used to distinguish AD and HC and are effective for AD detection.

## 5. Conclusions

In this paper, we expose an underlying bias present in the audio recordings produced from the picture description task of the Pitt corpus, a commonly used publicly available dataset for Alzheimer's Disease detection. Even by solely leveraging the silent segments of these audio recordings, we can achieve nearly 100% classification accuracy. Through experimental analysis, we propose that this bias is caused by background noise and other recording environment factors present in the original speech samples of the dataset. Subsequently, after preprocessing the speech samples with stationary noise removal and standard amplitude normalization, the experimental results demonstrate the alleviation of the bias, confirming the effectiveness of the data preprocessing methods. This study emphasizes the importance of understanding what the model has learned and calls for caution in the blind application of black-box automated models. We hope that researchers will pay attention to the potential dangers caused by spurious features in the data. In future work, we aim to delve into research on model interpretability, as well as explore cross-lingual Alzheimer's Disease detection utilizing both English and Mandarin datasets.

---

[5]https://github.com/jiaaro/pydub

# 6. References

[1] P. J. Nestor, P. Scheltens, and J. R. Hodges, "Advances in the early detection of alzheimer's disease," *Nature medicine*, vol. 10, no. Suppl 7, pp. S34–S41, 2004.

[2] J. T. Becker, F. Boiler, O. L. Lopez, J. Saxton, and K. L. McGonigle, "The natural history of alzheimer's disease: description of study cohort and accuracy of diagnosis," *Archives of neurology*, vol. 51, no. 6, pp. 585–594, 1994.

[3] H. Goodglass and E. Kaplan, *Boston diagnostic aphasia examination booklet*. Lea & Febiger, 1983.

[4] H. J. Han, S. BN, L. Qiu, and S. Abdullah, "Automatic classification of dementia using text and speech data," in *Multimodal AI in healthcare: A paradigm shift in health intelligence*. Springer, 2022, pp. 399–407.

[5] F. Chollet, "Xception: Deep learning with depthwise separable convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1251–1258.

[6] R. B. Ammar and Y. B. Ayed, "Evaluation of acoustic features for early diagnosis of alzheimer disease," in *Intelligent Systems Design and Applications: 19th International Conference on Intelligent Systems Design and Applications (ISDA 2019) held December 3-5, 2019 19*. Springer, 2021, pp. 172–181.

[7] S. Zargarbashi and B. Babaali, "A multi-modal feature embedding approach to diagnose alzheimer disease from spoken language," *arXiv preprint arXiv:1910.00330*, 2019.

[8] K. C. Fraser, J. A. Meltzer, and F. Rudzicz, "Linguistic features identify alzheimer's disease in narrative speech," *Journal of Alzheimer's Disease*, vol. 49, no. 2, pp. 407–422, 2016.

[9] S. Luz, F. Haider, S. de la Fuente, D. Fromm, and B. MacWhinney, "Alzheimer's dementia recognition through spontaneous speech: The adress challenge," *arXiv preprint arXiv:2004.06833*, 2020.

[10] S. Luz, F. Haider, S. de la Fuente, D. Fromm, and B. MacWhinney, "Detecting cognitive decline using speech only: The adresso challenge," *arXiv preprint arXiv:2104.09356*, 2021.

[11] K. Mei, Z. Guo, Z. Liu, L. Liu, X. Li, and Z. Ling, "Detecting alzheimer's disease based on acoustic features extracted from pretrained models," in *CAAI International Conference on Artificial Intelligence*. Springer, 2022, pp. 272–283.

[12] J. Koo, J. H. Lee, J. Pyo, Y. Jo, and K. Lee, "Exploiting multimodal features from pre-trained networks for alzheimer's dementia recognition," *arXiv preprint arXiv:2009.04070*, 2020.

[13] L. Gauder, L. Pepino, L. Ferrer, and P. Riera, "Alzheimer disease recognition using speech-based embeddings from pre-trained models." in *Interspeech*, 2021, pp. 3795–3799.

[14] J. Chen, J. Ye, F. Tang, and J. Zhou, "Automatic detection of alzheimer's disease using spontaneous speech only," in *Interspeech*, vol. 2021. NIH Public Access, 2021, p. 3830.

[15] O. Pfungst, *Clever Hans:(the horse of Mr. Von Osten.) a contribution to experimental animal and human psychology*. Holt, Rinehart and Winston, 1911.

[16] M. Arjovsky, L. Bottou, I. Gulrajani, and D. Lopez-Paz, "Invariant risk minimization," *arXiv preprint arXiv:1907.02893*, 2019.

[17] A. Borah, D. Pylypenko, C. Espana-Bonet, and J. van Genabith, "Measuring spurious correlation in classification:'clever hans' in translationese," *arXiv preprint arXiv:2308.13170*, 2023.

[18] B. Chettri, S. Mishra, B. L. Sturm, and E. Benetos, "Analysing the predictions of a cnn-based replay spoofing detection system," in *2018 IEEE spoken language technology workshop (SLT)*. IEEE, 2018, pp. 92–97.

[19] D. Wallis and I. Buvat, "Clever hans effect found in a widely used brain tumour mri dataset," *Medical Image Analysis*, vol. 77, p. 102368, 2022.

[20] C. Perlich, P. Melville, Y. Liu, G. Świrszcz, R. Lawrence, and S. Rosset, "Breast cancer identification: Kdd cup winner's report," *ACM SIGKDD Explorations Newsletter*, vol. 10, no. 2, pp. 39–42, 2008.

[21] A. J. DeGrave, J. D. Janizek, and S.-I. Lee, "Ai for radiographic covid-19 detection selects shortcuts over signal," *Nature Machine Intelligence*, vol. 3, no. 7, pp. 610–619, 2021.

[22] J. D. López-Cabrera, R. Orozco-Morales, J. A. Portal-Diaz, O. Lovelle-Enríquez, and M. Pérez-Díaz, "Current limitations to identify covid-19 using artificial intelligence with chest x-ray imaging," *Health and Technology*, vol. 11, no. 2, pp. 411–424, 2021.

[23] Z. Liu, Z. Guo, Z. Ling, S. Wang, L. Jin, and Y. Li, "Dementia detection by analyzing spontaneous mandarin speech," in *2019 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*. IEEE, 2019, pp. 289–296.

[24] Z. Guo, Z. Liu, Z. Ling, S. Wang, L. Jin, and Y. Li, "Text classification by contrastive learning and cross-lingual data augmentation for alzheimer's disease detection," in *Proceedings of the 28th international conference on computational linguistics*, 2020, pp. 6161–6171.

[25] Z. Sheng, Z. Guo, X. Li, Y. Li, and Z. Ling, "Dementia detection by fusing speech and eye-tracking representation," in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 6457–6461.

[26] H. Bredin, "pyannote.audio 2.1 speaker diarization pipeline: principle, benchmark, and recipe," in *Proc. INTERSPEECH 2023*, 2023.

[27] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," *Advances in neural information processing systems*, vol. 33, pp. 12 449–12 460, 2020.

[28] K. Mei, X. Ding, Y. Liu, Z. Guo, F. Xu, X. Li, T. Naren, J. Yuan, and Z. Ling, "The ustc system for adress-m challenge," in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023, pp. 1–2.

[29] F. Eyben, M. Wöllmer, and B. Schuller, "Opensmile: the munich versatile and fast open-source audio feature extractor," in *Proceedings of the 18th ACM international conference on Multimedia*, 2010, pp. 1459–1462.

[30] B. Schuller, S. Steidl, A. Batliner, J. Hirschberg, J. K. Burgoon, A. Baird, A. Elkins, Y. Zhang, E. Coutinho, and K. Evanini, "The interspeech 2016 computational paralinguistics challenge: Deception, sincerity & native language," in *17TH Annual Conference of the International Speech Communication Association (Interspeech 2016), Vols 1-5*, vol. 8. ISCA, 2016, pp. 2001–2005.

[31] Y. Jeon, J. Kang, B. C. Kim, K. H. Lee, J.-I. Song, and J. Gwak, "Early alzheimer's disease diagnosis using wearable sensors and multilevel gait assessment: A machine learning ensemble approach," *IEEE Sensors Journal*, 2023.

[32] T. Sainburg, M. Thielk, and T. Q. Gentner, "Finding, visualizing, and quantifying latent structure across diverse animal vocal repertoires," *PLoS computational biology*, vol. 16, no. 10, p. e1008228, 2020.