

# CHARACTERIZING THE ADVERSARIAL VULNERABILITY OF SPEECH SELF-SUPERVISED LEARNING

Haibin Wu<sup>12\*</sup>, Bo Zheng<sup>23\*</sup>, Xu Li<sup>3</sup>, Xixin Wu<sup>23</sup>, Hung-yi Lee<sup>1</sup>, Helen Meng<sup>23</sup>

<sup>1</sup> Graduate Institute of Communication Engineering, National Taiwan University

<sup>2</sup>Centre for Perceptual and Interactive Intelligence, The Chinese University of Hong Kong

<sup>3</sup>Human-Computer Communications Laboratory, The Chinese University of Hong Kong

## ABSTRACT

A leaderboard named Speech processing Universal Performance Benchmark (SUPERB), which aims at benchmarking the performance of a shared self-supervised learning (SSL) speech model across various downstream speech tasks with minimal modification of architectures and a small amount of data, has fueled the research for speech representation learning. The SUPERB demonstrates speech SSL upstream models improve the performance of various downstream tasks through just minimal adaptation. As the paradigm of the self-supervised learning upstream model followed by downstream tasks arouses more attention in the speech community, characterizing the adversarial robustness of such paradigm is of high priority. In this paper, we make the first attempt to investigate the adversarial vulnerability of such paradigm under the attacks from both zero-knowledge adversaries and limited-knowledge adversaries. The experimental results illustrate that the paradigm proposed by SUPERB is seriously vulnerable to limited-knowledge adversaries, and the attacks generated by zero-knowledge adversaries are with transferability. The XAB test verifies the imperceptibility of crafted adversarial attacks.

**Index Terms**— Adversarial attack, self-supervised learning

## 1. INTRODUCTION

There is a huge imbalance between labeled and unlabeled data, as labeled data is hard to obtain, yet unlabeled data is everywhere. Self-supervised learning (SSL) can take advantage of such large volumes of unlabeled data to mine general-purpose knowledge. It is a new trend in natural language processing [1] and computer vision [2] to pre-train a shared SSL upstream model followed by minimal adaptation to downstream tasks, and the features extracted by such upstream models will benefit the performance of downstream tasks. Recently, a leaderboard named Speech processing Universal Performance Benchmark (SUPERB) [3], which aims at benchmarking the performance of a shared speech self-supervised learning (SSL) model across various downstream speech tasks with minimal modification of architectures and small amount of data, has fueled the research for speech representation learning. Furthermore, SUPERB demonstrates that speech SSL upstream models can also improve and boost the performance of various downstream speech tasks through minimal adaptation. As the paradigm of the self-supervised learning upstream model followed by downstream tasks

\* Equal contribution. This research is funded by the Centre for Perceptual and Interactive Intelligence, an InnoCentre of The Chinese University of Hong Kong. This work was done while Haibin Wu was a visiting student at the CUHK. Haibin Wu is supported by Google PHD Fellowship Scholarship.

brings significant performance gains and arouses increasing attention in the speech community, it remains to be investigated whether the paradigm is robust enough to adversarial attacks.

The concept of adversarial attack was first proposed by [4], and the authors showed the state-of-the-art image classification models are vulnerable to adversarial attacks. Adversarial attacks are usually indistinguishable from their genuine counterparts based on human perception, yet it can manipulate the AI models and then cause them to have catastrophic failures. In this sense, adversarial attacks are particularly dangerous. Speech processing models, including automatic speech recognition (ASR) [5–8], automatic speaker verification (ASV) [9–16], anti-spoofing for ASV [17–21] and voice conversion [22], are also susceptible to adversarial attacks. Given a piece of audio, whether music, silence or speech, the authors [5] can generate adversarial audio, which is indistinguishable from the genuine version based on human’s ears, but will cause the ASR to transcribe any adversary-desired transcriptions. [9] and [13] respectively illustrate that adversarial attacks can also manipulate state-of-the-art ASV systems into falsely accepting the imposters or falsely rejecting the authorized persons. [18] is the first to show that the anti-spoofing model which shields ASV, are also vulnerable to adversarial attacks.

As SSL features attain the merits of generalizability and reusability, and the paradigm equipped with SSL achieves competitive performance in speech processing tasks, such a paradigm naturally arouses keen interests from both academia and industry. Whether such a paradigm would be an exception which can counter adversarial attacks remains an open question, and characterizing the adversarial robustness of such paradigm is of high priority. In this paper, we make the first attempt to investigate the adversarial vulnerability of such a paradigm under the attacks from both zero-knowledge adversaries and limited-knowledge adversaries (see section 3.1). The experiments mainly focus on attacking the upstream models (Sec 2.1), including wav2vec 2.0 [23] and HuBERT [24], without access to the downstream models (Sec 2.3) and task-specific procedures (Sec 2.2), in order to issue the attack across tasks. The experimental results show the paradigm proposed by SUPERB is vulnerable to limited-knowledge adversaries, and the attacks generated by zero-knowledge adversaries are transferable. The XAB test verifies the imperceptibility of the carefully concocted adversarial attacks.

## 2. UPSTREAM-DOWNSTREAM PARADIGM

SUPERB first introduced the upstream-downstream paradigm for speech processing in a systematic view. The paradigm is shown in Fig. 1 (b). The self-supervised learning models which learn general-purpose knowledge from a large amount of unlabeled data play the role of upstream models (Sec 2.1). The upstream models are pre-

trained and then the parameters are frozen during downstream models training and inference. A task-specific module (Sec 2.2), consisting of a layer-wise weighted sum procedure and a pre-processing procedure, is designed for each downstream task. The downstream models (Sec 2.3) get the features  $z$  and  $\tilde{z}$ , rather than traditional acoustic features, e.g. MFCC.

## 2.1. Upstream

### 2.1.1. HuBERT

HuBERT adopts BERT-style token classification for pre-training. Off-line unsupervised clustering is first applied to acoustic features, such as MFCC, to get frame-level noise labels. Then the extracted features from convolutional layers are masked, and based on the noise labels, BERT-like predictive loss is applied to the masked regions. The authors expect the pre-trained models can create better features than MFCC, so they re-implement the clustering to the HuBERT features in the early training iterations to get better noise labels, then repeat the BERT-like pre-training. HuBERT gets the best performance in the SUPERB.

### 2.1.2. Wav2vec 2.0

wav2vec 2.0 learns general-purpose knowledge by contrastive loss. It firstly masks the hidden speech representations extracted by a multi-layer convolutional network from an utterance, followed by transformer layers to build contextualized representations given the hidden representations. After quantization of the hidden representations to derive the latent for each hidden representations, a contrastive task is introduced to distinguish the true latent and the distractors. wav2vec 2.0 achieves comparable performance with that of HuBERT across the SUPERB downstream tasks.

## 2.2. Task-specific module

The task-specific module is composed of a layer-wise weighted sum procedure and a pre-processing procedure. The layer-wise weighted sum procedure consists of task-specific weights applied to all the hidden features from different layers of the upstream model, and such weights are usually jointly trained with downstream models. The pre-processing procedure is designed for each downstream task, such as pre-emphasis and voice activity detection. As we proceed to training and testing for the downstream task, the audio data first undergo task-specific pre-processing, and are then fed into the upstream model to extract hidden features, followed by the layer-wise weighted sum to get the final features - these are then adopted to train the downstream model. However, during adversarial attack, the adversaries generally do not pre-process the audio, and only use averaged embeddings from each layer of the upstream model to make the attack less task-specific.

## 2.3. Downstream Tasks

**Phoneme Recognition (PR)** recognizes phonemes from an utterance. SUPERB selects train-clean-100, dev-clean, and test-clean subsets of LibriSpeech [25] as training, validation, and testing set, respectively. Phone error rate (PER) is the evaluation metric.

**Automatic Speech Recognition (ASR)** recognizes words from an utterance. SUPERB adopts train-clean-100, dev-clean, and test-clean subsets of LibriSpeech [25] as training, validation, and testing set, respectively. Word error rate (WER) is the evaluation metric.

**Keyword Spotting (KS)** identifies predefined keywords from an utterance. SUPERB uses Speech Commands dataset v1.0 [26], which introduces 12 classes, including 10 for keywords, one for silence, and one for unknown words. Accuracy (ACC) is the evaluation metric.

**Speaker Identification (SID)** is a close-set multi-class classification task, which aims at identifying the speaker of a given utterance from a set of predefined speakers. VoxCeleb1 [27] is adopted by SUPERB. Accuracy (ACC) is the evaluation metric.

**Automatic Speaker Verification (ASV)** verifies whether a pair of utterances belong to the same speaker. VoxCeleb1 [27] is adopted in this task. Accuracy (ACC), rather than equal error rate (EER) is the evaluation metric during attack (refer to 4.1).

**Speaker Diarization (SD)** identifies who is speaking during each timestamp in an utterance. LibriMix [28] is used in this task, which is generated from LibriSpeech [25]. The speaker labels of each chunk are generated by alignments from Kaldi [29] LibriSpeech ASR model. Diarization error rate (DER) is the evaluation metric.

**Intent Classification (IC)** identifies predefined classes of intent from an utterance. Fluent Speech Commands [30] is used, where every utterance is labeled with one of the three intent classes: action, object, and location. Accuracy (ACC) is the evaluation metric.

**Slot Filling (SF)** identifies all semantic slots of an utterance with predefined slot-type and slot-value. Audio SNIPS [31] is adopted by SUPERB to generate multi-speaker utterances from SNIPS [32]. According to the standard split in SNIPS, samples of US-accent speakers are selected as training set, and others as validation and testing sets. Since slot-type and slot-value are both essential in SF task, F1 score and CER [33] are adopted as evaluation metrics for slot-type and slot-value, respectively.

**Emotion Recognition (ER)** recognizes emotion class from an utterance. IEOMCAP [34] is used. Following the past evaluation protocol, SUPERB discards unbalanced classes, resulting in four emotion classes remaining: neutral, happy, sad, and angry. The evaluation metric is accuracy (ACC).

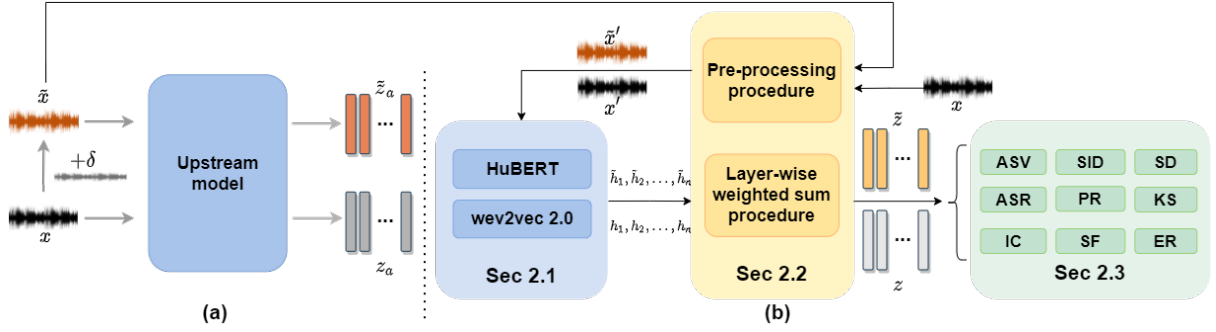
The models of downstream tasks are simply structured in SUPERB. A linear network is used for PR, and optimized by CTC loss. For KS, SID, IC, and ER, linear transformation models after mean-pooling trained by cross-entropy loss are adopted. As for ASR, a 2-layer 1024-unit BLSTM is used, and optimized by CTC loss. For SF, slot-type labels are represented into an ordered pair wrapping the corresponding slot-value to form a sequence of tokens, and then SF is transformed into an ASR task using the same model as ASR. For ASV, x-vector [35] is the backbone model optimized with AM-Softmax loss [36], and cosine-similarity backend is used for scoring. For SD, a single-layer 512-unit LSTM is used with permutation-invariant training (PIT) loss [37].

## 3. ADVERSARIAL ATTACK FOR SPEECH SSL

### 3.1. Attacking scenarios

In this work, we distinguish different attack scenarios from the perspective of the knowledge accessed by adversaries.

**Limited-knowledge adversaries:** Attackers can access the internals of the target upstream model, including the detailed parameters and gradients. But they do not know which downstream task will be conducted, not to mention the internals of downstream models, the weights of the layer-wise weighted sum procedure and task-specific pre-processing procedures. Knowing the internals of the target model, the attackers will directly calculate the gradients and generate adversarial samples, as shown in Fig. 1.(a).



**Fig. 1.** (a) Attacking framework for SSL.  $x$  and  $\tilde{x}$  are the original and adversarial samples respectively,  $z_a$  and  $\tilde{z}_a$  are the features of the upstream model given  $x$  and  $\tilde{x}$  as inputs respectively, and  $\delta$  is the carefully designed adversarial noise. (b) Upstream-downstream paradigm.  $x$  and  $x'$  are the original and pre-processed samples respectively,  $h_1, h_2, \dots, h_n$  are hidden features extracted from upstream model, where the subscript of  $h_n$  denotes the layer number of the upstream model, and  $z$  is the final features obtained from weighted sum of  $h_1, h_2, \dots, h_n$  by the layer-wise weighted sum procedure.  $\tilde{x}, \tilde{x}', \tilde{h}_1, \tilde{h}_2, \dots, \tilde{h}_n$  and  $\tilde{z}$  are the adversarial counterparts of  $x, x', h_1, h_2, \dots, h_n$  and  $z$  respectively.

**Zero-knowledge adversaries:** In this scenario, while the attackers aim at the *target model*, it is unavailable to the attackers. In such a case, the *substitute model* is used for approximating gradients for adversarial sample generation. Zero-knowledge adversaries can get even less knowledge than limited-knowledge adversaries. Zero-knowledge adversaries do not even know the details of the target upstream model internals. In order to conduct adversarial attacks, they have to train another substitute model, adopt the gradients of the substitute model to generate adversarial samples, and finally use such adversarial samples to fool the upstream-downstream paradigm using the target model.

### 3.2. Attack procedure

Under the scenarios of zero-knowledge attacks and limited-knowledge attacks, attackers only have access to the upstream model without knowing the task-specific module (Section 2.2) and the downstream model (Section 2.3) to craft the adversarial attacks. For limited-knowledge adversaries, they have access to the target upstream model internals, while the zero-knowledge adversaries will train a substitute upstream model and use it to generate adversarial attacks.

During an attack, the weights in layer-wise weighted sum procedure are set as equal to average the embeddings in each layer of the upstream model to derive  $z_a$  and  $\tilde{z}_a$ . Whatever the downstream task is, we only manipulate the upstream model to generate adversarial samples. So in such a scenario, the attack method introduced below is less task-specific. The attacking framework is shown in Fig. 1(a). Having fixed parameters in the upstream model, the attackers aim at crafting the adversarial noise  $\delta$  to maximize the difference between  $z_a$  and  $\tilde{z}_a$ , while also let  $x$  and  $\tilde{x}$  indistinguishable from human's ears. In order to fulfill the above two objectives, we introduce the basic iterative method (BIM) [38] for attack. BIM crafts the adversarial sample in an iterative manner. Starting from the genuine input  $x^0 = x$ , the input is perturbed iteratively as

$$x^{n+1} = \text{clip}_{x,\epsilon}(x^n + \delta), \quad (1)$$

for  $n = 0, \dots, N - 1$

where  $\delta$  is the carefully designed adversarial noise, derived as

$$\delta = \alpha \times \text{sign}(\nabla_{x^n} \|z_a - \tilde{z}_a\|_2) \quad (2)$$

where  $\alpha$  is the step size,  $\text{sign}(\cdot)$  is a function which gets the sign of the gradient,  $N$  is the number of iterations and  $\text{clip}_{x,\epsilon}(t)$  is the

norm constraint which conducts element-wise clipping such that  $\|t - x\|_\infty \leq \epsilon$  to assure the original sample  $x$  and the derived adversarial sample  $x^N$  are indistinguishable.  $x^N$  will be used for attacking the upstream-downstream paradigm.

## 4. EXPERIMENT

### 4.1. Experimental setup

We train the upstream models in Fig. 1.(b), and fix its parameters as the downstream models are trained. We omit the implementation details due to space limitation, and readers can refer to [3] for more information. Note that the adversarial attacks are conducted during inference. As the adversarial attack is time- and resource-consuming, we randomly selected 100 genuine samples for attacking, do the experiments three times, and then report the mean and variance of results. Note that for ASV, we randomly select 50 non-target and target trials. The ACC for ASV is derived by the number of trials with the right decision over the total trial number. The performance for the genuine samples is shown in rows (d) and (h) of Table 1. The standard deviations are less than 0.1%, so we only show the means in rows (d) and (h). Then we craft adversarial samples according to the attacking methods as in Section 3.2. Gaussian noise of the same noise-to-signal ratio (NSR) with adversarial perturbations is introduced for comparison (in rows (c) and (g)).

### 4.2. Experimental Results

Table 1 illustrates the attack performance on SSL for 9 downstream tasks. The direction of the arrow in the second row denotes the direction towards the better performance of the task. For example,  $\downarrow$  for ASR means the lower WER implies better ASR performance, yet the less effective the attacks are. The first column in Table 1 lists the *method to generate the attack model and the target model*. For example, w2v2-w2v2 denotes that the substitute model for generating adversarial samples is wav2vec 2.0, and the target model is also wav2vec 2.0. The rows (a) and (e) are the limited-knowledge scenarios. The rows (b) and (f) are the zero-knowledge scenarios. gau-w2v2 denotes using the samples perturbed by Gaussian noise to attack wav2vec 2.0. We adopt gau-w2v2 and gau-HuBERT as our baseline, and the results are in rows (c) and (g).

**Table 1.** Adversarial attack performance on SSL representations for various downstream tasks.

		ASR	PR	KS	IC	SF		SID	ER	SD		ASV
		WER ↓	PER ↓	Acc ↑	Acc ↑	F1 ↑	CER ↓	Acc ↑	Acc ↑	Acc ↑	DER ↓	Acc ↑
(a)	w2v2-w2v2	19.20 <sup>1</sup> (±2.01)	28.32 (±2.03)	65.67 (±6.51)	55.67 (±5.77)	88.55 (±1.33)	20.19 (±2.05)	81.33 (±3.06)	79.33 (±3.79)	88.48 (±0.19)	17.48 (±0.55)	91.67 (±2.31)
(b)	HuBERT-w2v2	5.54 (±0.71)	5.09 (±0.47)	91.00 (±3.00)	88.33 (±1.15)	95.36 (±1.26)	8.70 (±0.55)	87.67 (±4.16)	87.33 (±6.03)	94.56 (±0.36)	8.08 (±0.41)	97.00 (±2.00)
(c)	gau-w2v2	0.48 (±0.06)	1.11 (±0.05)	98.67 (±0.58)	93.67 (±1.15)	99.71 (±0.27)	0.71 (±0.50)	97.67 (±2.08)	95.67 (±3.06)	98.24 (±0.09)	2.51 (±0.11)	99 (±0.00)
(d)	Clean-w2v2	0	0	100	100	100	0	100	100	100	2.51	100
(e)	HuBERT-HuBERT	26.76 (±0.82)	18.67 (±1.54)	64.33 (±0.58)	69.67 (±5.03)	76.91 (±1.67)	36.54 (±1.83)	76.33 (±4.93)	78.33 (±2.08)	87.78 (±0.83)	18.39 (±1.65)	88.33 (±2.08)
(f)	w2v2-HuBERT	1.94 (±0.06)	2.21 (±0.28)	96.67 (±1.15)	98.33 (±1.15)	99.42 (±0.37)	1.62 (±0.16)	93.67 (±1.15)	91.00 (±2.65)	95.13 (±0.20)	7.17 (±0.47)	96.67 (±1.53)
(g)	gau-HuBERT	0.05 (±0.08)	0.42 (±0.12)	99.67 (±0.58)	99.67 (±0.58)	99.89 (±0.19)	0.25 (±0.24)	98.67 (±2.31)	99.00 (±0.00)	98.36 (±0.09)	2.32 (±0.13)	99.67 (±0.58)
(h)	Clean-HuBERT	0	0	100	100	100	0	100	100	98.37	2.31	100

<sup>1</sup> We show the mean and standard deviation. Here 19.20 ±2.01 means that the mean and standard deviation of WER are 19.20% and 2.01%.

**Table 2.** Adversarial attack performance ASR.

	NSR	EDR	WER
w2v2-w2v2	0.67(±0.01)	81.15(±0.27)	19.20(±2.01)
HuBERT-w2v2	0.71(±0.00)	58.09(±0.09)	5.54(±0.71)
gau-w2v2	0.69(±0.00)	30.40(±0.19)	0.48(±0.06)

We have these observations: (1) Limited-knowledge attackers achieve the most effective attack on wav2vec 2.0 and HuBERT for all downstream tasks. Take the IC task as an example, “w2v2-w2v2” degrades the Acc of wav2vec 2.0 from 100% to 55.67%, and “HuBERT-HuBERT” degrades that of HuBERT from 100% to 69.67%. These results verify the severe threats that adversarial attacks can pose on the SSL models. We also observe similar trends for all other 8 tasks. (2) Zero-knowledge attackers achieve relatively weaker attacks on downstream tasks than limited-knowledge attackers, but the attack is still effective. Specifically, in some downstream tasks, zero-knowledge attackers also seriously degrade the system performance. For instance, in the ER downstream task, “HuBERT-w2v2” degrades the Acc of wav2vec 2.0 from 100% to 87.33%, and “w2v2-HuBERT” degrades the Acc of HuBERT from 100% to 91.00%. (3) To verify the effectiveness of adversarial perturbations, we leverage the Gaussian noise for comparison. We observe that Gaussian noise degrades the wav2vec 2.0 and HuBERT much less for all downstream tasks compared with both limited-knowledge and zero-knowledge attacks. For instance, in the ASR task, Gaussian noise has little effect on the wav2vec 2.0 WER degradation (0.48%) and practically has no influence on the HuBERT WER degradation (0.05%). This suggests that simply adding Gaussian noise cannot degrade a well-trained system for the malicious attack purpose and verifies the effectiveness and transferability of our attacks.

Moreover, we also compare the NSR and embedding distance rate (EDR) for the limited-knowledge attackers, zero-knowledge attackers and Gaussian noise in Table 2. EDR is calculated by  $1/N \times \sum_{n=1}^N \|z_n - \tilde{z}_n\|_2 / \|z_n\|_2$ , where  $z_n$  and  $\tilde{z}_n$  are the adversarial-original embedding pair,  $N$  is the total adversarial-original pair number. Here we only show the results on the ASR task with wav2vec 2.0 as target model due to space limitation, while other tasks have similar trends. Table 2 illustrates the results. We set the NSR of all perturbations at a similar scale for fair comparison. We observe that the EDR

has similar trends with the WER. Specifically, limited-knowledge attackers make the embeddings of adversarial inputs most far away from the original ones, resulting in the largest EDR when compared with the zero-knowledge attackers and Gaussian noise. For instance, “w2v2-w2v2” achieves an EDR of 81.15%. For zero-knowledge attackers, “HuBERT-w2v2” achieves an EDR of 58.09%, which is less than those of the limited-knowledge attackers. While for Gaussian noise, “gau-w2v2” only achieves an EDR of 30.40%. This suggests the less effectiveness of Gaussian noise. Finally, we observe a consistency between EDR and the system performance WER.

### 4.3. XAB test

To illustrate the imperceptibility of adversarial noise, we conduct the XAB listening test. XAB listening test is a standard test to evaluate the detectability between two choices of sensory stimuli. We randomly select 2 adversarial-genuine pairs for each upstream-downstream paradigm. 36 randomly selected adversarial-genuine pairs (i.e., A and B) are shown to the listeners. We randomly choose one from each pair as the reference audio (i.e., X) and let the listeners select the audio, which sounds more similar to X, from A and B. Five listeners take part in the XAB listening test. The XAB test has a classification accuracy of 58.89%, which illustrates that the adversarial samples are hard to be distinguished from genuine samples. Audio demos with the attacking settings are made open here <sup>2</sup>.

## 5. CONCLUSION

Adversarial robustness is an AI standard for trustworthy machine learning systems. Though the paradigm proposed by SUPERB is going to penetrate all the speech processing tasks and gains good performance, the adversarial robustness has not received sufficient consideration. This work is the first to expose the vulnerability of such paradigm to adversarial attacks. In future work, we will investigate attacks with higher transferability and less imperceptibility. As more sophisticated attacks continue to be developed, we will need to come up with defense methods to alleviate such attacks. The long-term goal is to design adaptive defense methods that offer protection against increasingly dangerous attacks.

<sup>2</sup>Audio demo

## 6. REFERENCES

- [1] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.
- [2] A. Newell and J. Deng, "How useful is self-supervised pre-training for visual tasks?," in *Proceedings of the IEEE/CVF CVPR*, 2020, pp. 7345–7354.
- [3] S.-w. Yang et al., "Superb: Speech processing universal performance benchmark," *arXiv preprint arXiv:2105.01051*, 2021.
- [4] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," *arXiv preprint arXiv:1312.6199*, 2013.
- [5] N. Carlini and D. Wagner, "Audio adversarial examples: Targeted attacks on speech-to-text," in *2018 IEEE Security and Privacy Workshops (SPW)*. IEEE, 2018, pp. 1–7.
- [6] H. Yakura and J. Sakuma, "Robust audio adversarial example for a physical attack," *arXiv preprint arXiv:1810.11793*, 2018.
- [7] R. Taori, A. Kamsetty, B. Chu, and N. Vemuri, "Targeted adversarial examples for black box audio systems," in *SPW 2019*. IEEE, 2019, pp. 15–20.
- [8] Y. Qin, N. Carlini, G. Cottrell, I. Goodfellow, and C. Raffel, "Imperceptible, robust, and targeted adversarial examples for automatic speech recognition," in *International Conference on Machine Learning*. PMLR, 2019, pp. 5231–5240.
- [9] J. Villalba, Y. Zhang, and N. Dehak, "x-vectors meet adversarial attacks: Benchmarking adversarial robustness in speaker verification," *Proc. Interspeech 2020*, pp. 4233–4237, 2020.
- [10] F. Kreuk, Y. Adi, M. Cisse, and J. Keshet, "Fooling end-to-end speaker verification with adversarial examples," in *ICASSP*. IEEE, 2018, pp. 1962–1966.
- [11] M. Marras, P. Korus, N. D. Memon, and G. Fenu, "Adversarial optimization for dictionary attacks on speaker verification," in *Interspeech*, 2019, pp. 2913–2917.
- [12] H. Wu et al., "Improving the adversarial robustness for speaker verification by self-supervised learning," *IEEE/ACM TASLP*, vol. 30, pp. 202–217, 2021.
- [13] X. Li et al., "Adversarial attacks on gmm i-vector based speaker verification systems," in *ICASSP*. IEEE, 2020, pp. 6579–6583.
- [14] H. Wu et al., "Voting for the right answer: Adversarial defense for speaker verification," in *Interspeech*, 2021.
- [15] H. Wu et al., "Adversarial defense for automatic speaker verification by cascaded self-supervised learning models," in *ICASSP*. IEEE, 2021, pp. 6718–6722.
- [16] H. Wu et al., "Spotting adversarial samples for speaker verification by neural vocoders," *arXiv preprint arXiv:2107.00309*, 2021.
- [17] A. Kassis and U. Hengartner, "Practical attacks on voice spoofing countermeasures," *arXiv preprint arXiv:2107.14642*, 2021.
- [18] S. Liu, H. Wu, H.-y. Lee, and H. Meng, "Adversarial attacks on spoofing countermeasures of automatic speaker verification," in *2019 IEEE ASRU*. IEEE, 2019, pp. 312–319.
- [19] H. Wu, S. Liu, H. Meng, and H.-y. Lee, "Defense against adversarial attacks on spoofing countermeasures of asv," *arXiv preprint arXiv:2003.03065*, 2020.
- [20] H. Wu, A. T. Liu, and H.-y. Lee, "Defense for black-box attacks on anti-spoofing models by self-supervised learning," in *Interspeech*, 2020.
- [21] Y. Zhang, Z. Jiang, J. Villalba, and N. Dehak, "Black-box attacks on spoofing countermeasures using transferability of adversarial examples," *Interspeech 2020*, pp. 4238–4242, 2020.
- [22] C.-y. Huang, Y. Y. Lin, H.-y. Lee, and L.-s. Lee, "Defending your voice: Adversarial attack on voice conversion," in *SLT 2021*. IEEE, 2021, pp. 552–559.
- [23] A. Baevski, H. Zhou, A. Mohamed, and M. Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," *arXiv preprint arXiv:2006.11477*, 2020.
- [24] W.-N. Hsu, Y.-H. H. Tsai, B. Bolte, R. Salakhutdinov, and A. Mohamed, "Hubert: How much can a bad teacher benefit asr pre-training?," in *ICASSP*. IEEE, 2021, pp. 6533–6537.
- [25] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: An asr corpus based on public domain audio books," *ICASSP 2015*, pp. 5206–5210, 2015.
- [26] P. Warden, "Speech commands: A dataset for limited-vocabulary speech recognition," *arXiv preprint arXiv:1804.03209*, 2018.
- [27] A. Nagrani, J. S. Chung, and A. Zisserman, "Voxceleb: a large-scale speaker identification dataset," *arXiv preprint arXiv:1706.08612*, 2017.
- [28] J. Cosentino, M. Pariente, S. Cornell, A. Deleforge, and E. Vincent, "Librimix: An open-source dataset for generalizable speech separation," *arXiv preprint arXiv:2005.11262*, 2020.
- [29] D. Povey et al., "The kaldi speech recognition toolkit," in *ASRU*. IEEE Signal Processing Society, 2011, number CONF.
- [30] L. Lugosch, M. Ravanelli, P. Ignoto, V. S. Tomar, and Y. Bengio, "Speech model pre-training for end-to-end spoken language understanding," *arXiv preprint arXiv:1904.03670*, 2019.
- [31] C.-I. Lai, Y.-S. Chuang, H.-Y. Lee, S.-W. Li, and J. Glass, "Semi-supervised spoken language understanding via self-supervised speech and language model pretraining," *arXiv preprint arXiv:2010.13826*, 2020.
- [32] A. S. Alice Coucke et al., "Snips voice platform: an embedded spoken language understanding system for private-by-design voice interfaces," *arXiv preprint arXiv:1805.10190*, 2018.
- [33] N. Tomashenko et al., "Recent advances in end-to-end spoken language understanding," *Lecture Notes in Computer Science*, p. 44–55, 2019.
- [34] M. B. Carlos Busso et al., "Iemocap: interactive emotional dyadic motion capture database," *Language Resources and Evaluation*, vol. 42, pp. 335–359, 2008.
- [35] D. Snyder, D. Garcia-Romero, G. Sell, D. Povey, and S. Khudanpur, "X-vectors: Robust dnn embeddings for speaker recognition," in *ICASSP 2018*. IEEE, 2018, pp. 5329–5333.
- [36] H. Wang et al., "Cosface: Large margin cosine loss for deep face recognition," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 5265–5274.
- [37] Y. Fujita et al., "End-to-end neural speaker diarization with permutation-free objectives," *arXiv preprint arXiv:1909.05952*, 2019.
- [38] A. Kurakin, I. Goodfellow, and S. Bengio, "Adversarial examples in the physical world," *arXiv preprint arXiv:1607.02533*, 2016.