# FAST YET EFFECTIVE SPEECH EMOTION RECOGNITION WITH SELF-DISTILLATION

*Zhao Ren*[1], *Thanh Tam Nguyen*[2], *Yi Chang*[3], *Björn W. Schuller*[3,4]

[1]L3S Research Center, Leibniz University Hannover, Germany
[2]Griffith University, Australia
[3]GLAM – Group on Language, Audio, & Music, Imperial College London, United Kingdom
[4]Chair of Embedded Intelligence for Health Care and Wellbeing, University of Augsburg, Germany

`zren@l3s.de`

## ABSTRACT

Speech emotion recognition (SER) is the task of recognising human's emotional states from speech. SER is extremely prevalent in helping dialogue systems to truly understand our emotions and become a trustworthy human conversational partner. Due to the lengthy nature of speech, SER also suffers from the lack of abundant labelled data for powerful models like deep neural networks. Pretrained complex models on large-scale speech datasets have been successfully applied to SER via transfer learning. However, fine-tuning complex models still requires large memory space and results in low inference efficiency. In this paper, we argue achieving a fast yet effective SER is possible with self-distillation, a method of simultaneously fine-tuning a pretrained model and training shallower versions of itself. The benefits of our self-distillation framework are threefold: (1) the adoption of self-distillation method upon the acoustic modality breaks through the limited ground-truth of speech data, and outperforms the existing models' performance on an SER dataset; (2) executing powerful models at different depth can achieve adaptive accuracy-efficiency trade-offs on resource-limited edge devices; (3) a new fine-tuning process rather than training from scratch for self-distillation leads to faster learning time and the state-of-the-art accuracy on data with small quantities of label information.

***Index Terms***— self-distillation, speech emotion recognition, adaptive inference, efficient deep learning, efficient edge analytics

## 1. INTRODUCTION

Speech emotion recognition (SER) nowadays is an idiosyncratic task in many dialogue systems, such as Siri, Cortana, and Alexa [1]. Through classifying human speech signals into various emotional states (e. g., happiness, surprise, anger, disgust, fear, sadness, neutral, etc.), SER helps human-computer systems become more personalised and trustworthy as well as adjust the contexts accordingly in car-driving, heath-diagnosis, call-center, aircraft-cockpit, and web/mobile applications [2, 3].

Existing techniques for SER are limited by the inherent lack of labelled data due to the expensive efforts of annotation (e.g. thousands of hours of speech over nearly 7,000 spoken languages [4]). They often rely on large deep neural networks that are pre-trained by

unsupervised learning, contrastive learning, or self-supervised learning, such as wav2vec [5], wav2vec 2.0 [4], and vq-wav2vec [6]. However, fine-tuning large models has a high demand of memory space and inference time [7].

In machine learning, self-distillation has emerged as a paradigm to develop a student model with a more lightweight architecture that can even outperform the teacher [8]. This has been particularly successfully applied to computer vision [8, 9]. However, in contrast to the visual modality, the acoustic modality is significantly more challenging due to limited ground-truth. Self-distillation methods cannot be applied directly to SER since they often require large labelled data to simultaneously train a teacher model from scratch with shallower student versions of itself [8].

In this paper, we present a framework of self-distillation for fast, yet effective speech emotion recognition. While our framework is demonstrated on wav2vec 2.0 [4] (one of the state-of-the-art (SOTA) pre-trained models for speech representations), it can be applied to other models and datasets with limited ground-truth information. In our framework (see Figure 1), the pre-trained wav2vec 2.0 (i. e., the teacher model) was fine-tuned together with shallower model parameters from itself (i. e., the student models), when the teacher and all students are predicting emotional states from speech samples.

To the best of our knowledge, this is the first attempt to develop a self-distillation framework for SER. The contributions of our self-distillation framework include: (1) the application of self-distillation on speech data overcomes the difficulty caused by limited annotations, and outperforms the existing models' performance on an SER dataset; (2) executing powerful models at different depths increases the possibility to achieve adaptive accuracy-efficiency trade-offs on resource-limited edge devices; (3) a new fine-tuning process rather than training from scratch for self-distillation leads to faster learning time and SOTA accuracy on data with limited ground-truth.

**Related Works.** Spectrum features [10, 11] have been often used as the input of deep neural networks for SER [12], while selecting the appropriate spectrum features is a time-consuming work. Moreover, the performance of SER is limited to expensive human annotations; lacking of labelled data for deep learning. More recently, self-supervised learning on speech data has shown promising to learn effective representations, and the pre-trained models have been successfully fine-tuned for SER tasks [13–15]. Therefore, we apply an end-to-end self-supervised learning model, wav2vec 2.0, to SER.

Knowledge distillation is one of the popular methods to achieve high efficiency by transferring knowledge from a teacher model to a smaller student model [7]. Similar to other model compression approaches such as pruning and quantisation, they sacrifice information loss (thus accuracy) and could not overcome the accuracy-

efficiency trade-offs. Our self-distillation approach can achieve the best of both worlds by reusing the architecture and allowing inference at different depths of the teacher model itself.

Different types of self-distillation methods have been developed recently, including iteration-based [16, 17], aggregate-based [18], and branch-based approaches [8, 9]. Iteration-based methods perform knowledge distillation from a teacher model to a student model with the same architecture and this procedure is repeated a few times [16, 17]. However, the training and inference costs are not reduced, as the teacher and the student are the same. Aggregated-based methods use data augmentation to produce more versions of the teacher model on different augmentations and then combine the outputs [18]. However, existing data augmentations are domain-specific and are not applicable to acoustic modality. Our work relates closely to the branch-based approaches, which add branches at different depths of the teacher model using bottlenecks/attention-blocks and shallow classifiers [8, 9]. However, these layers are not applicable for wave2vec 2.0, as it already contains transformer layers. Our work is also different from the layer-wise knowledge distillation, which fine-tunes a student model from the teacher model itself for predicting deep layers of the teacher [7]. The layer-wise knowledge distillation can produce a general student model, while it requires fine-tuning efforts for a specific task. The fixed number of model parameters of the student is limited for performance improvement with a deeper structure and lacks of flexibility.

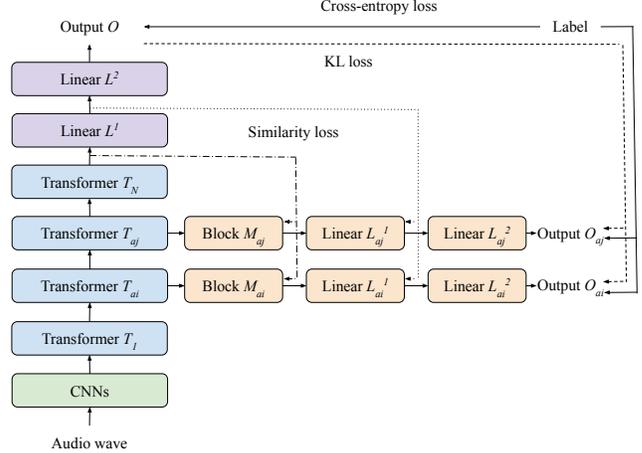## 2. METHODOLOGY

### 2.1. Preliminaries

**Self-supervised Learning with wav2vec 2.0.** Self-supervised learning has shown its superiority compared to supervised learning on many audio tasks, such as speaker recognition [19] and SER [20]. Wav2vec 2.0 [4] was trained on the large-scale Librispeech corpus [21] in a self-supervised learning framework. Wav2vec 2.0 has been widely used to extract effective representations for SER tasks [13, 22]. A Wav2vec 2.0 model consist of multi-layer convolutional neural networks (CNNs) (i. e., encoder) and multiple transformer layers (i. e., context network). The latent representations learnt from the encoder are discretised into a set of quatisation representations, which are then processed with the output of the context network in a contrastive task.

### 2.2. Self-Distillation Framework: The Case of Wav2vec 2.0

#### 2.2.1. Model Architecture

**Teacher Model.** We assume the input data of wav2vec 2.0 is represented as $(\boldsymbol{X}, y)$, where $\boldsymbol{X}$ is the raw speech signals and $y$ is the emotional states. The final (i. e., $N$-th) transformer layer $T_N$ of wav2vec 2.0 is followed by two linear layers ($L^1$ and $L^2$) with output dimensions of $D_1$ and $D_2$, where $D_2$ is the number of emotional classes (see Figure 1). Regarding the intermediate features, the output of each transformer layer has a dimension of $(B, F, R)$, where $B$ is the sample number in a batch, $F$ represents the number of time steps, and $R$ denotes the dimension of representations at each time step. With the goal of classification, the $N$-th transformer layer's output is pooled into $H_N$ with a dimension of $(B, R)$ before being fed into the two linear layers.

**Student Model.** To reduce the model parameters of wav2vec 2.0 with self-distillation, additional layers are added after the intermediate transformer layers of wav2vec 2.0 (see Figure 1). In a student



**Fig. 1**: The framework of self-distillation on wav2vec 2.0. The output of an intermediate transformer layer is processed by a block model and two linear layers. The output of each second linear layer is the predicted probabilities on emotional classes. The backward process is implemented via three types of loss functions. '—': cross-entropy loss, '- - -': Kullback–Leibler (KL) divergence loss, '- · - ·' or '· · · · ·': similarity loss.

model, the transformer layer $T_{ai}$, $1 \leq ai < N$, is followed by a distillation model, including a block $M_{ai}$ and two linear layers ($L^1_{ai}$ and $L^2_{ai}$). $M_{ai}$ is a neural network, and $L^2_{ai}$ is trained for predicting emotional classes. Herein, as $M_{ai}$ is expected to learn representations similar to those from $T_N$, the output of $T_{ai}$ is directly fed into the block model without pooling. Therefore, the output of $M_{ai}$ has a dimension of $(B, F, R)$, and is pooled into $H_{M_{ai}}$ with a dimension of $(B, R)$. $H_{M_{ai}}$ is then fed into $L^1_{ai}$ for further process.

Apart from the above single distillation model, multiple distillation models could be learnt together in self-distillation to improve the flexibility for different depths of models. For instance, two distillation models in Figure 1 are developed after transformer layers $T_{ai}$ and $T_{aj}$ to build SER models with different numbers of model parameters.

#### 2.2.2. Loss Function

As the pre-trained wav2vec 2.0 has already strong capability of learning representations from speech, we initialise the teacher model's parameters with the pre-trained wav2vec 2.0. We assume $N_a$ distillation models are built after transformer layers indexed by $A = \{a1, a2, ...\}$. The representations $H_L$ and $H_{L_{ai}}$ are learnt from $L^1$ and $L^1_{ai}$, respectively. The outputs of the second linear layers $L^2$ and $L^2_{ai}$ are represented as $O$ and $O_{ai}$. The model parameters are optimised with the loss function in self-distillation:

$$\mathcal{L} = \mathcal{L}_c + \alpha \mathcal{L}_k + \beta \mathcal{L}_s, \qquad (1)$$

where $\mathcal{L}_c$ is the cross-entropy loss, $\mathcal{L}_k$ is the Kullback–Leibler (KL) divergence loss, $\mathcal{L}_s$ is the similarity loss, and $\alpha$ and $\beta$ are constant values.

**Cross-entropy Loss.** To train wav2vec 2.0 and distillation models for performing SER, the cross-entropy loss contains i) the cross entropy loss on the teacher model (i. e., $L^1$, $L^2$, and wav2vec 2.0), and ii) the cross entropy loss on the student models (i. e., distillation models and partial model parameters of wav2vec 2.0):

$$\mathcal{L}_c = \mathcal{L}_{\text{ce}}(O, y) + \gamma \frac{1}{N_a} \sum_{ai \in A} \mathcal{L}_{\text{ce}}(O_{ai}, y), \qquad (2)$$

where $\mathcal{L}_{\text{ce}}$ is the typical cross entropy loss for training a model in supervised learning, and $\gamma$ is a constant value.

**KL Loss.** The outputs of the distillation models are expected to be similar to that of the linear layer $L^2$, which is the final layer of the teacher model. With this target, the Kullback-Leibler (KL) loss aims to regularise the outputs $O$ and $O_{ai}$:

$$\mathcal{L}_k = \frac{1}{N_a} \sum_{ai \in A} O \log \frac{O}{O_{ai}}. \tag{3}$$

**Similarity Loss.** Apart from the loss functions computed on the model outputs, loss functions on the interval features learnt from the intermediate layers can further help train strong student models. In this work, we compare three loss functions, including $L_1$, $L_2$, and cosine similarity. Their combinations are also compared with single functions. Furthermore, these loss functions could be either on the outputs of the $N$-th transformer layer and the blocks ($H_N$ and $H_{M_{ai}}$), or on the output of the first linear layers ($H_L$ and $H_{L_{ai}}$):

$$\mathcal{L}_s = \frac{1}{N_a} \sum_{ai \in A} \mathcal{L}_{\text{sim}}(H_N, H_{M_{ai}}) \text{ OR } \frac{1}{N_a} \sum_{ai \in A} \mathcal{L}_{\text{sim}}(H_L, H_{L_{ai}}), \tag{4}$$

where $\mathcal{L}_{\text{sim}}$ is a loss function of $L_1$, $L_2$ or negative cosine similarity.

### 2.3. Dynamic Inference

Although the deep layers of wav2vec 2.0 can often learn higher-level reresentations than shallow ones, self-distillation can provide dynamic inference models [9] via training both shallow and deep student models with good performance. Shallow distillation models require less parameters than deep ones, and deep ones may perform better than shallow ones. The flexibility of self-distillation enables SER applications to be applied to various hardwares, from wearable devices to work stations.
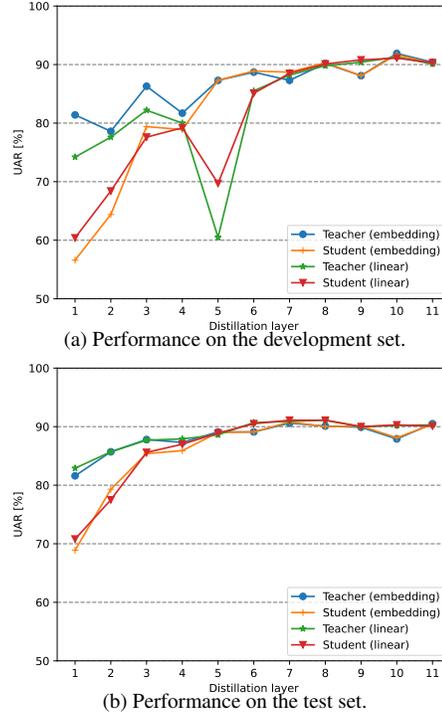
## 3. EXPERIMENTS

### 3.1. Database

The database of elicited mood in speech (DEMoS) [23] is used to verify the self-distillation for SER. DEMoS with 9 365 emotional and 332 neutral Italian speech samples was collected from 68 speakers (f: 23, m: 45) [23]. Each speech sample was annotated with one of the eight classes: *anger*, *disgust*, *fear*, *guilt*, *happiness*, *sadness*, *surprise*, and *neutral*. To implement experiments that can be compared to other studies [24, 25] on DEMoS, the minor *neutral* class is not used, and the speaker-independent training/development/test sets are the same to our prior study [24]. The detail of the data distribution on the seven emotional classes can be found in [24].

### 3.2. Experimental Settings

**Evaluations Metrics.** In this study, the unweighted average recall (UAR) is employed to evaluate the performance of SER models. A UAR is computed as the average of all class-wise recalls.

**Implementation Details.** With the goal of classifying emotion states, the wav2vec 2.0 model is followed by two linear layers with the numbers of output neurons $\{256, 7\}$, respectively. In each distillation model, the block could be one of the three layers: CNN (number of output channels: 1, kernel size: $(1, 1)$), Long Short-Time Memory recurrent neural networks (LSTM-RNN) (number of output



(a) Performance on the development set.

(b) Performance on the test set.

**Fig. 2**: Comparison of the performance (UAR [%]) of the teacher models and the student ones, each of which has single convolutional distillation layer. The loss function is $L_2$ loss.

features: 768), Gated Recurrent Unit (GRU) RNN (number of output features: 768). The two linear layers in each distillation model also have the numbers of output neurons $\{256, 7\}$, respectively.

Each model is trained on the training set and validated on the development set, and further trained on the combination of the training and development sets, and validated on the test set. During training, the hyperparameters of the loss function are set as $\alpha = \beta = \gamma = 1$. All training procedures of self-distillation are optimised by an Adam optimiser with a learning rate of $3e - 5$, and stopped at the 20-th epoch, when the batch size is 16.

**Reproducibility Environments.** To improve the reproducibility, the code of this work is released at: `https://github.com/leibniz-future-lab/SelfDistill-SER`.

### 3.3. Sensitivity Analysis

Figure 2 shows the performance of self-distillation with single CNN-based distillation model. We can see that, on both development set (Figure 2 (a)) and test set (Figure 2 (b)), the performance is increasing when the distillation layer is going deeper. This indicates that deeper model layers of wav2vec 2.0 can learn more abstract representations than shallower ones. The teacher models perform better than shallow student models, but are comparable with deep student models, especially after the 7-th distillation layer.

As the distillation layers at the embedding level perform slightly better than those at the linear level, the embedding level is selected in the next experiments for self-distillation. To provide different model sizes with self-distillation, we group the distillation layers into three groups and select layers which have the best performance on the development set. Therefore, we use layers $\{3, 8, 10\}$ for multi-layer self-distillation in the following experiments.

**Table 1**: Comparison of the performance (UAR [%]) of multi-layer self-distillation (distillation layers: {3, 8, 10}). The best performance of the deepest model is highlighted with '‗', and the best performance of distillation layers is highlighted with bold fonts.

| NN | Loss | Deepest | | Layer 3 | | Layer 8 | | Layer 10 | | Fusion | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Devel | Test | Devel | Test | Devel | Test | Devel | Test | Devel | Test |
| Self-distillation w/ CNN | $L_1$ | 87.3 | 89.0 | 75.2 | 82.4 | 86.9 | 88.8 | 87.1 | 88.9 | 85.7 | 88.3 |
| | $L_2$ | 88.3 | 89.3 | 70.2 | 81.6 | 88.6 | 89.4 | 88.2 | 89.2 | 86.8 | 88.5 |
| | Cosine sim. | 89.6 | 89.7 | **76.0** | **81.5** | **89.5** | **89.9** | **89.7** | **89.7** | **88.7** | **89.8** |
| | $L_1$ + Cosine sim. | 88.4 | 87.8 | 77.4 | 77.5 | 88.4 | 88.1 | 88.4 | 87.8 | 87.5 | 87.4 |
| | $L_2$ + Cosine sim. | 88.1 | 89.7 | 77.5 | 81.5 | 88.8 | 89.3 | 88.4 | 89.5 | 87.4 | 89.0 |
| Self-distillation w/ LSTM | $L_1$ | 91.4 | 90.2 | **83.3** | **85.7** | **91.7** | **90.7** | **91.8** | **90.5** | **90.5** | **89.9** |
| | $L_2$ | 90.7 | 91.2 | 78.5 | 86.7 | 91.0 | 91.1 | 90.8 | 90.9 | 89.6 | 90.8 |
| | Cosine sim. | 90.8 | 88.8 | 79.2 | 75.2 | 91.1 | 87.9 | 90.8 | 88.5 | 90.1 | 87.1 |
| | $L_1$ + Cosine sim. | 90.3 | 90.6 | 77.7 | 86.5 | 90.1 | 90.3 | 90.3 | 90.5 | 88.8 | 90.5 |
| | $L_2$ + Cosine sim. | 91.8 | 91.4 | 83.3 | 84.7 | 91.1 | 91.4 | 91.7 | 91.2 | 90.1 | 90.9 |
| Self-distillation w/ GRU | $L_1$ | 86.3 | 91.0 | 80.2 | 85.5 | 84.7 | 91.6 | 84.5 | 91.2 | 86.1 | 90.8 |
| | $L_2$ | 90.9 | 90.7 | **82.0** | **84.1** | **91.2** | **90.4** | **91.1** | **90.6** | **90.0** | **90.0** |
| | Cosine sim. | 90.2 | 89.6 | 74.8 | 82.2 | 90.3 | 89.5 | 89.9 | 89.4 | 88.7 | 88.9 |
| | $L_1$ + Cosine sim. | 91.2 | 90.0 | 79.5 | 83.7 | 91.4 | 90.2 | 91.0 | 90.1 | 89.5 | 89.7 |
| | $L_2$ + Cosine sim. | 88.1 | 91.6 | 72.2 | 85.1 | 87.5 | 91.5 | 86.5 | 91.7 | 86.8 | 90.7 |

## 3.4. Ablation Study

With the distillation layers {3, 8, 10}, we compare the three block models (i. e., CNN, LSTM-RNN, and GRU-RNN) with various similarity loss functions (i. e., $\mathcal{L}_{sim}$). From Table 2, we can find that all similarity loss functions perform comparably for each block model. Particularly, the single similarity loss functions (i. e., $L_1$, $L_2$, and cosine similarity) perform better than the combinations of them. This might be related to the setting of hyperparameters in the loss functions. Furthermore, the LSTM-RNN and the GRU-RNN models outperform the CNN one when comparing the three models blocks. This may be because RNNs can better learn sequential information than CNNs. Regarding the self-distillation, the performance of the student models at layers 8 and 10 is comparable with the deepest model. Layer 3 performs slightly worse than layers 8 and 10, as the corresponding student model of layer 3 is shallower than those of layers 8 and 10. Finally, the fusion results of the three distillation models are comparable with the results of layer 10.

## 3.5. Comparison with SOTA

We compared the results of self distillation with the other SOTA models, including the following three groups of models. (1) CNN-4, VGG-16, ResNet-50, and VGG-16 with adversarial training are trained from scratch [24]. (2) The models of wav2vec 2.0 with fine-tuning are trained for 20 epochs based on the pre-trained wav2vec 2.0, when the transformer layers after finetuned layers are frozen. (3) The layer-wise models are trained via the layer-wise knowledge distillation in [7]. As wav2vec 2.0 has 12 transformer layers, which is the same as the number of encoder layers in HuBERT in [7], the layer-wise models are mostly developed according to the settings in [7]. The teacher model is the pre-trained wav2vec 2.0, and the student model is part of the pre-trained wav2vec 2.0 itself (from the first layer to the second transformer layer). Notably, all parameters of the teacher model are frozen. The layer-wise models are trained at two stages: i) training the student model to predict layers {4, 8, 12}) of the teacher model, and ii) fine-tuning the student model on DEMoS for SER. To train a strong student model, the first stage is trained with 20 epochs. To implement fair experimental comparisons, the second stage also consists of 20 epochs.

As wav2vec 2.0 was pre-trained on the large-scale speech database, The models based on wav2vec 2.0 are mostly better than models trained from scratch (CNN-4, VGG-16 (+ adversarial training), and ResNet-50). When comparing fine-tuned wav2vec 2.0

**Table 2**: Comparison of the performance (UAR [%]) between our approach (lower lines) and the state-of-the-art (SOTA).

| NN | Devel | Test | #Param |
|---|---|---|---|
| CNN-4 [24] | 82.6 | 83.6 | 4.3 M |
| VGG-16 [24] | 79.8 | 83.6 | 14.7 M |
| ResNet-50 [24] | 71.9 | 81.3 | 23.5 M |
| VGG-16 + adversarial training [24] | 87.5 | 86.7 | 14.7 M |
| Wav2vec2 (layer 3) + fine-tuning | 77.1 | 83.4 | 31.5 M |
| Wav2vec2 (layer 8) + fine-tuning | 90.1 | 90.9 | 66.9 M |
| Wav2vec2 (layer 10) + fine-tuning | 91.7 | 89.2 | 81.1 M |
| Wav2vec2 (deepest) + fine-tuning | 91.1 | 90.6 | 95.2 M |
| Layer-wise distillation w/ CNN | 54.3 | 72.5 | 24.4 M |
| Layer-wise distillation w/ LSTM | 70.8 | 79.0 | 38.5 M |
| Layer-wise distillation w/ GRU | 73.1 | 77.4 | 35.0 M |
| **Self-distillation (layer 3)** | 83.3 | 85.7 | 36.2 M |
| **Self-distillation (layer 8)** | 91.7 | 90.7 | 71.6 M |
| **Self-distillation (layer 10)** | 91.8 | 90.5 | 85.8 M |
| **Self-distillation (teacher)** | **91.8** | **91.4** | 100.0 M |

models and self-distillation, the student model at layer 3 outperforms the corresponding fine-tuned one. The fine-tuned models are comparable with self-distillation at layers 8 and 10, while self-distillation trains student models at different layers in single training procedure. Our self-distillation outperforms layer-wise distillation at all three layers. This may be caused by the shallower student models (encoder and two transformer layers) in layer-wise distillation. For a specific task, self-distillation requires less training epochs (20 in our work) than layer-wise distillation (40 in our study), increasing the training efficiency.

## 4. CONCLUSIONS AND FUTURE WORK

This work aimed to reduce model parameters via self-distillation for fast and effective speech emotion recognition. The experiments were implemented on the Database of Elicited Mood in Speech (DE-MoS) [23] with the pre-trained wav2vec 2.0. The experimental results demonstrated that the student model at a shallow layer (layer 3) outperformed the corresponding fine-tuned wav2vec 2.0, and self-distillation achieved comparable performance with that of fine-tuned wav2vec 2.0 at deep layers. Moreover, self-distillation performed better than layer-wise knowledge distillation. In future work, the self-distillation approach will be verified on multiple databases. We will also investigate to further reduce the wav2vec 2.0 model by the state-of-the-art model compression approaches [26, 27].

# 5. REFERENCES

[1] Björn W Schuller, "Speech emotion recognition: Two decades in a nutshell, benchmarks, and ongoing trends," *Communications of the ACM*, vol. 61, no. 5, pp. 90–99, 2018.

[2] Moataz El Ayadi, Mohamed S Kamel, and Fakhri Karray, "Survey on speech emotion recognition: Features, classification schemes, and databases," *Pattern recognition*, vol. 44, no. 3, pp. 572–587, 2011.

[3] Ruhul Amin Khalil, Edward Jones, Mohammad Inayatullah Babar, Tariqullah Jan, Mohammad Haseeb Zafar, and Thamer Alhussain, "Speech emotion recognition using deep learning techniques: A review," *IEEE Access*, vol. 7, pp. 117327–117345, 2019.

[4] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," in *Proc. NeurIPS*, Vancouver, Canada, 2020, pp. 1–12.

[5] Steffen Schneider, Alexei Baevski, Ronan Collobert, and Michael Auli, "wav2vec: Unsupervised pre-training for speech recognition," in *Proc. INTERSPEECH*, Graz, Austria, 2019, pp. 3465–3469.

[6] Alexei Baevski, Steffen Schneider, and Michael Auli, "vq-wav2vec: Self-supervised learning of discrete speech representations," in *ICLR*, Virtual, 2020.

[7] Heng-Jui Chang, Shu-wen Yang, and Hung-yi Lee, "Distilhubert: Speech representation learning by layer-wise distillation of hidden-unit bert," in *Proc. ICASSP*, Singapore, 2022, pp. 7087–7091.

[8] Linfeng Zhang, Jiebo Song, Anni Gao, Jingwei Chen, Chenglong Bao, and Kaisheng Ma, "Be your own teacher: Improve the performance of convolutional neural networks via self distillation," in *Proc. ICCV*, Seoul, Korea, 2019, pp. 3713–3722.

[9] Linfeng Zhang, Chenglong Bao, and Kaisheng Ma, "Self-distillation: Towards efficient and compact neural networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 8, pp. 4388–4403, 2021.

[10] Yuan Gong, Cheng-I Lai, Yu-An Chung, and James Glass, "Ssast: Self-supervised audio spectrogram transformer," in *Proc. AAAI*, Virtual, 2022, vol. 36, pp. 10699–10709.

[11] Prabhav Singh, Ridam Srivastava, KPS Rana, and Vineet Kumar, "A multimodal hierarchical approach to speech emotion recognition from audio and text," *Knowledge-Based Systems*, vol. 229, pp. 107316, 2021.

[12] Taiba Majid Wani, Teddy Surya Gunawan, Syed Asif Ahmad Qadri, Mira Kartiwi, and Eliathamby Ambikairajah, "A comprehensive review of speech emotion recognition systems," *IEEE Access*, vol. 9, pp. 47795–47814, 2021.

[13] Leonardo Pepino, Pablo Riera, and Luciana Ferrer, "Emotion recognition from speech using wav2vec 2.0 embeddings," in *Proc. INTERSPEECH*, Brno, Czech Republic, 2021, pp. 3400–3404.

[14] Li-Wei Chen and Alexander Rudnicky, "Exploring wav2vec 2.0 fine-tuning for improved speech emotion recognition," *arXiv preprint arXiv:2110.06309*, 2021.

[15] Johannes Wagner, Andreas Triantafyllopoulos, Hagen Wierstorf, Maximilian Schmitt, Florian Eyben, and Björn W Schuller, "Dawn of the transformer era in speech emotion recognition: closing the valence gap," *arXiv preprint arXiv:2203.07378*, 2022.

[16] Hossein Mobahi, Mehrdad Farajtabar, and Peter Bartlett, "Self-distillation amplifies regularization in hilbert space," in *Proc. NeurIPS*, Vancouver, Canada, 2020, pp. 1–11.

[17] Minh Pham, Minsu Cho, Ameya Joshi, and Chinmay Hegde, "Revisiting self-distillation," *arXiv preprint arXiv:2206.08491*, 2022.

[18] Hankook Lee, Sung Ju Hwang, and Jinwoo Shin, "Self-supervised label augmentation via input transformations," in *Proc. ICML*, Virtual, 2020, pp. 5714–5724.

[19] Nik Vaessen and David A Van Leeuwen, "Fine-tuning wav2vec2 for speaker recognition," in *Proc. ICASSP*, Singapore, 2022, pp. 7967–7971.

[20] Omar Mohamed and Salah A Aly, "Arabic speech emotion recognition employing wav2vec2. 0 and hubert based on baved dataset," *arXiv preprint arXiv:2110.04425*, 2021.

[21] Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur, "Librispeech: An ASR corpus based on public domain audio books," in *Proc. ICASSP*, Brisbane, Australia, 2015, pp. 5206–5210.

[22] Konlakorn Wongpatikaseree, Sattaya Singkul, Narit Hnoohom, and Sumeth Yuenyong, "Real-time end-to-end speech emotion recognition with cross-domain adaptation," *Big Data and Cognitive Computing*, vol. 6, no. 3, pp. 79, 2022.

[23] Emilia Parada-Cabaleiro, Giovanni Costantini, Anton Batliner, Maximilian Schmitt, and Björn Schuller, "DEMoS: An Italian emotional speech corpus," *Language Resources and Evaluation*, pp. 1–43, Feb. 2019.

[24] Zhao Ren, Alice Baird, Jing Han, Zixing Zhang, and Björn Schuller, "Generating and protecting against adversarial attacks for deep speech-based emotion recognition models," in *Proc. ICASSP*, Barcelona, Spain, 2020, pp. 7184–7188.

[25] Zhao Ren, Jing Han, Nicholas Cummins, and Björn Schuller, "Enhancing transferability of black-box adversarial attacks via lifelong learning for speech emotion recognition models," in *Proc. INTERSPEECH*, Shanghai, China, 2020, pp. 496–500.

[26] Anthony Berthelier, Thierry Chateau, Stefan Duffner, Christophe Garcia, and Christophe Blanc, "Deep model compression and architecture optimization for embedded systems: A survey," *Journal of Signal Processing Systems*, vol. 93, no. 8, pp. 863–878, 2021.

[27] Ke Tan and DeLiang Wang, "Towards model compression for deep learning based speech enhancement," *IEEE/ACM transactions on audio, speech, and language processing*, vol. 29, pp. 1785–1794, 2021.