



The Importance of Calibration: Rethinking Confidence and Performance of Speech Multi-label Emotion Classifiers

Huang-Cheng Chou^{1,2}, Lucas Goncalves¹, Seong-Gyun Leem¹,
Chi-Chun Lee², Carlos Busso¹

¹ Multimodal Signal Processing (MSP) lab,
Department of Electrical and Computer Engineering, University of Texas at Dallas (UTD), USA

² Behavioral Informatics & Interaction Computation (BIIC) lab,
Department of Electrical Engineering, National Tsing Hua University (NTHU), Taiwan

What is model calibration?

- It is a method to modify the predictions of models to improve the consistence between model accuracies and predictions' probabilities
 - If the classifier is well-calibrated, given 100 predictions, each with a confidence of 0.8, we expect that 80% of them should be correctly classified

Why does model calibration matter?

- Guo et al. [1] discovered that predictions of modern neural networks are often over-confident in computer vision and document classification tasks
 - e.g., high confidence for predictions with low accuracies

When does model calibration matter?

- Access probabilities of predictions for a richer interpretation
 - e.g., analyze the model shortcomings or provide the uncertainty to the end-users

[1] Guo et al.. (2017, July). On calibration of modern neural networks. In *International conference on machine learning* (pp. 1321-1330). PMLR.

Do modern Speech Emotion Classifiers (SECs) need model calibrations?:

- Yes. Our preliminary results show that **SECs** using one [1] of the state-of-the-art frameworks are under-confident
- Reliability Scenario: the most intuitive way to improve the reliability of predictions is to reject predictions [2,3] based on probabilities
- Issue: we may reject too many “low confident” predictions for samples that are actually correctly predicted

Hypotheses: the following factors could improve SECs’ calibration and classification performances

- Consider emotion co-occurrence
- Deal with the imbalance of emotional classes
- A multi-label post-calibration *temperature scaling* (TS) method

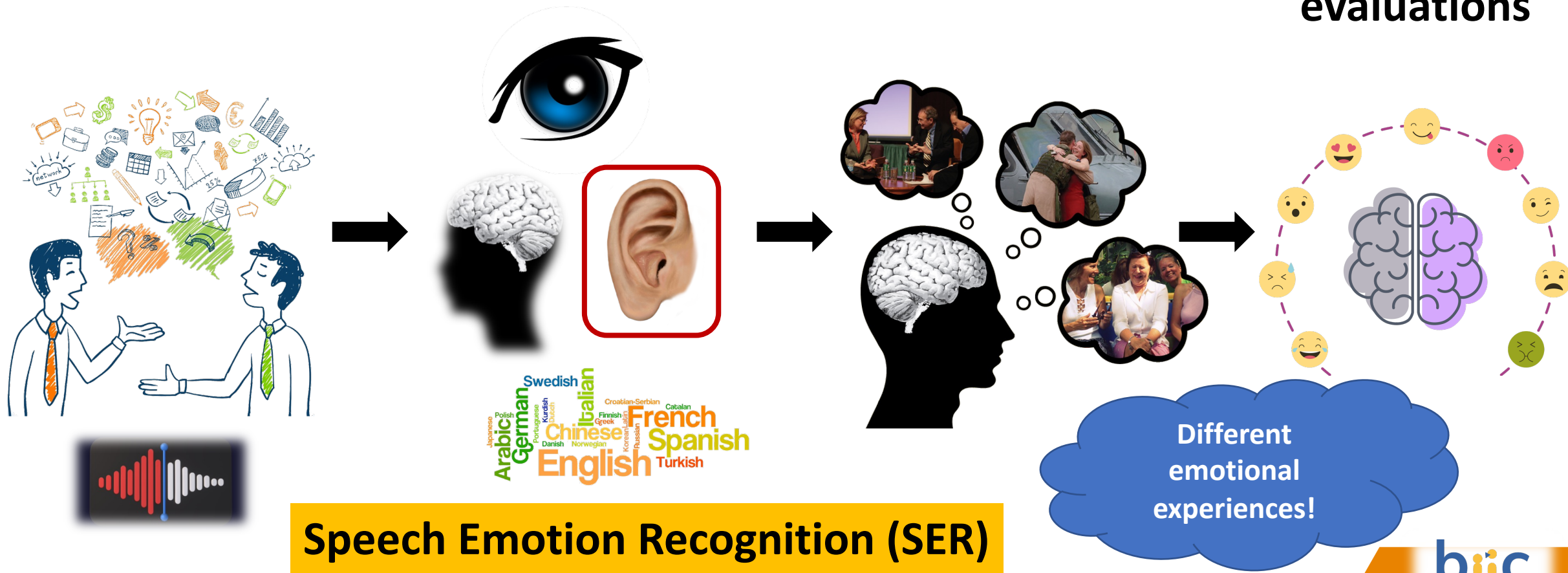
[1] Wagner et al. (2023). Dawn of the transformer era in speech emotion recognition: closing the valence gap. IEEE Transactions on Pattern Analysis and Machine Intelligence.

[2] Sridhar, K., & Busso, C. (2019, September). Speech Emotion Recognition with a Reject Option. In *Interspeech* (pp. 3272-3276).

[3] Sridhar, K., & Busso, C. (2020, May). Modeling uncertainty in predicting emotional attributes from spontaneous speech. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (pp. 8384-8388). IEEE.

Emotion Co-occurrence - Subjectivity of Emotion Perception

Emotional stimulus Emotion perception Emotion decoding Perceptual evaluations



Task: 8-class Primary Emotion:

File name: 0004_0073.wav

Annotations:

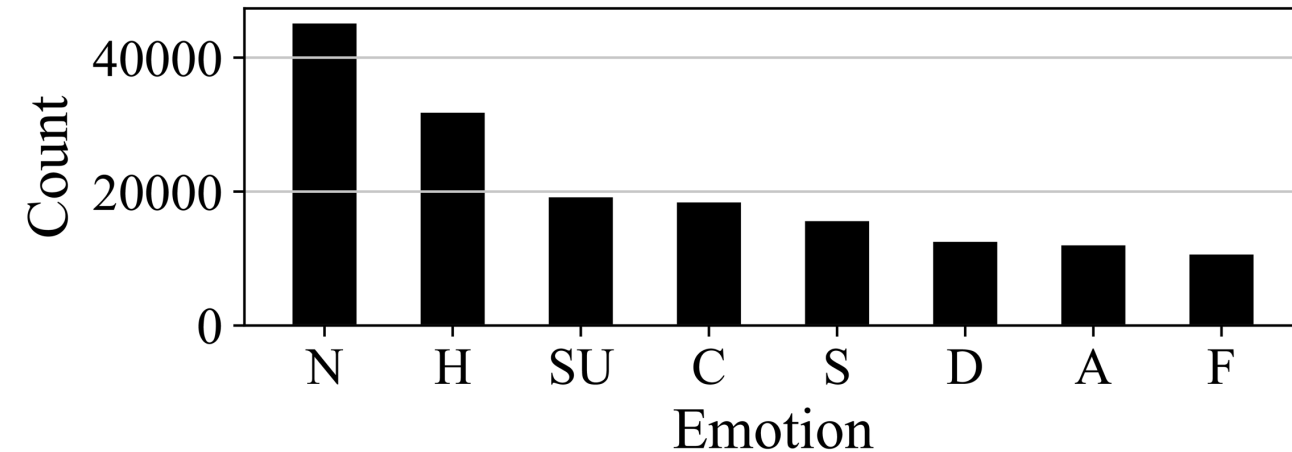
- Rater 1: Neutral
- Rater 2: Neutral
- Rater 3: Sad
- Rater 4: Angry
- Rater 5: Disgust

Quantity of classes on one sample:

- Neutral: 2
 - Sad: 1
 - Angry: 1
 - Disgust: 1
- **Emotion co-occurrence**
 - **Disagreements among raters**
 - **One sample has multiple annotations**

Imbalance annotation distribution:

- Zhong et al. [1] found that the class imbalance makes a model more miscalibrated
- The right-hand side figure shows the imbalance in the emotional distribution of the MSP-Podcast corpus based on the multi-label setting
 - One sentence might have more than one emotional label
- Most previous studies on SER have ignored imbalanced class distribution



N: neutral

S: sadness

H: happiness

D: disgust

SU: surprise

A: anger

C: contempt

F: fear

- **Temperature scaling calibration [1] is the most common post-calibration way to calibrate models**
- It was originally used for single-label tasks
- We adapt it for multi-label tasks, and we will introduce in detail in the next sessions

Purpose:

- Explore whether **considering emotion co-occurrence, imbalance of emotional annotations, and a calibration method can make models better calibrated and improve performance**

Method:

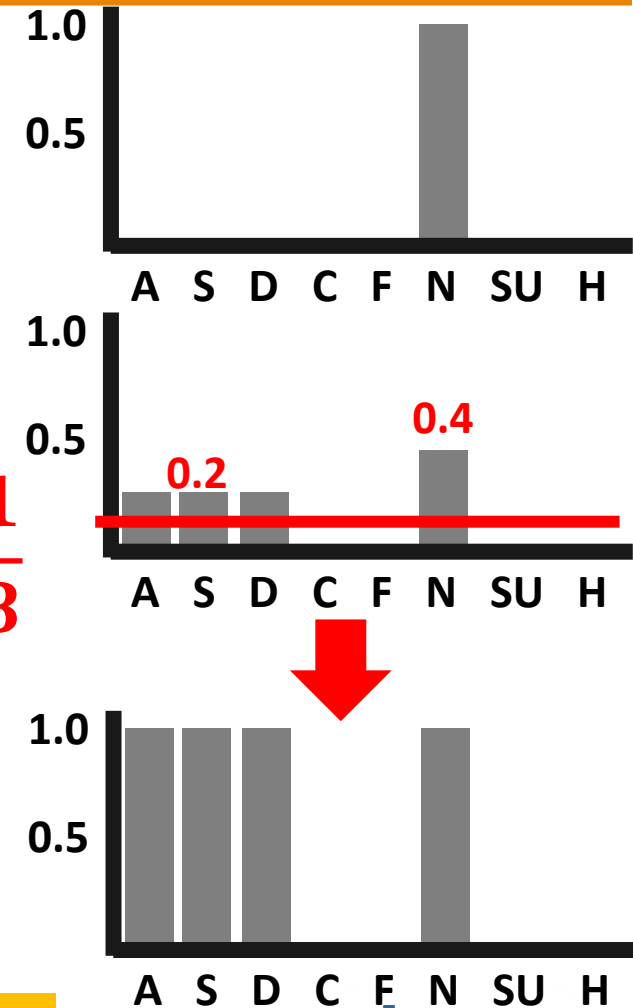
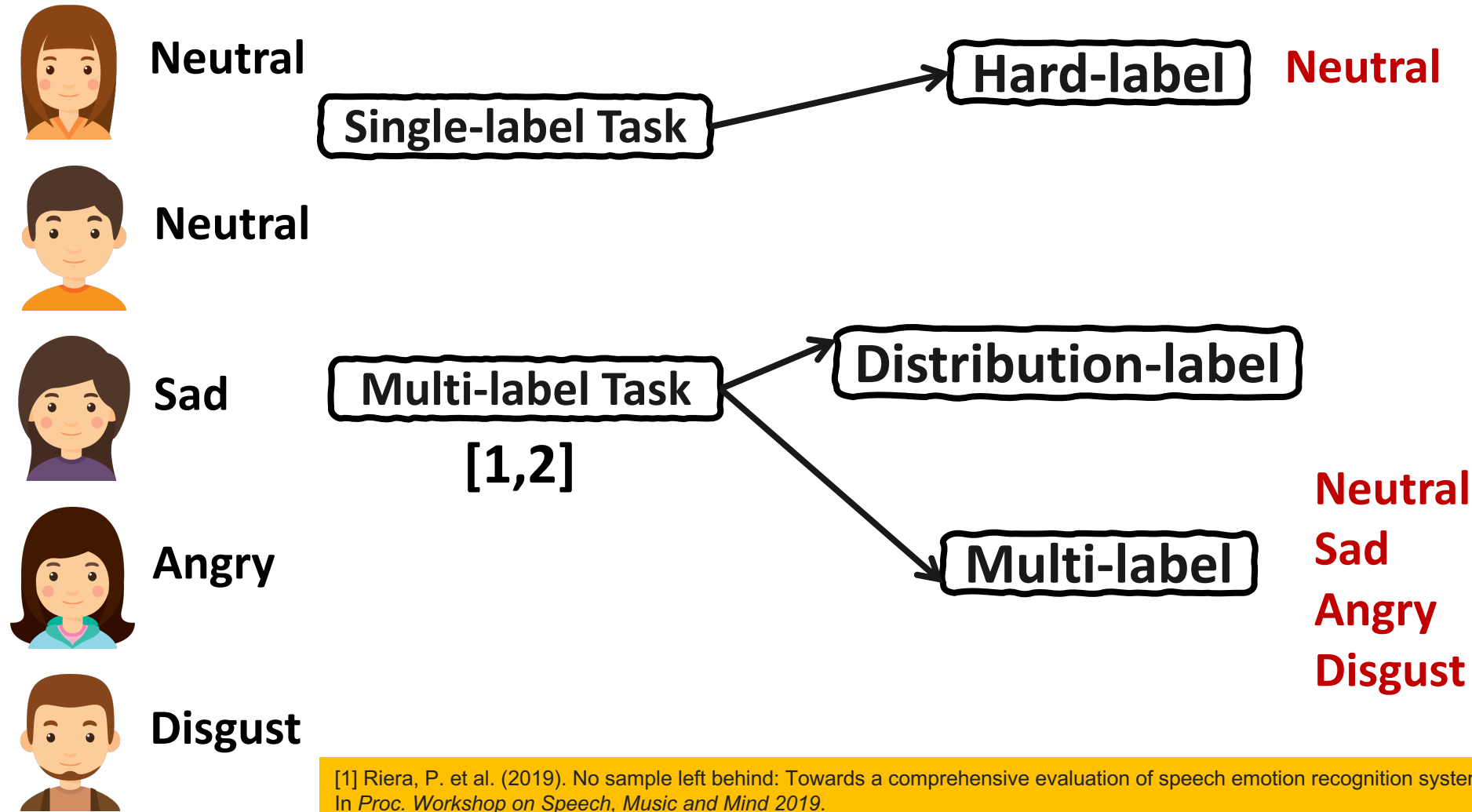
- **Emotion co-occurrence:** jointly training with emotion co-occurrence weight penalty loss function [1]
- **Imbalance of emotional annotations:** jointly training with class-balanced loss function [2]
- **Post-Calibration:** modifying the existing *temperature scaling* (TS) calibration [3] **for multi-label classification tasks**

[1] Chou, H. C., Lee, C. C., & Busso, C. (2022). Exploiting co-occurrence frequency of emotions in perceptual evaluations to train a speech emotion classifier. In *Proc. Interspeech* (Vol. 2022)..

[2] Cui, et al. (2019). Class-balanced loss based on effective number of samples. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 9268-9277).

[3] Guo et al.. (2017, July). On calibration of modern neural networks. In *International conference on machine learning* (pp. 1321-1330). PMLR.

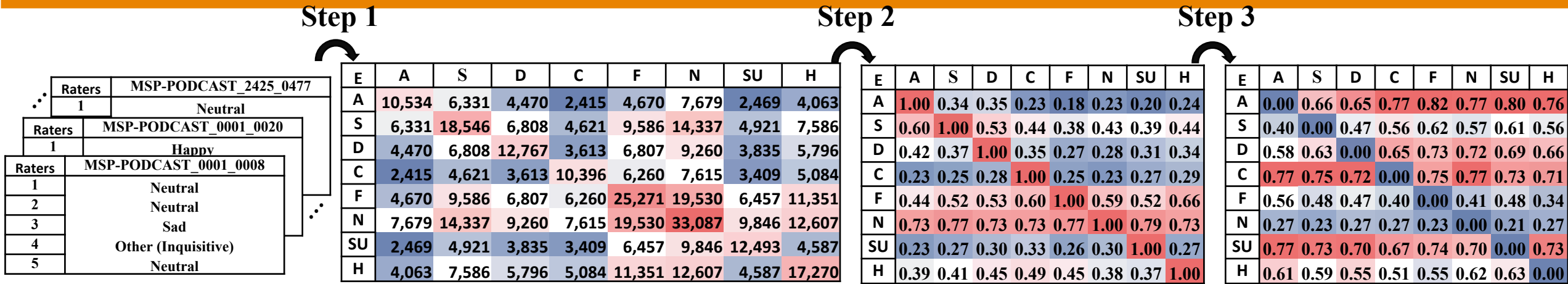
Decision of Labels for Speech Emotion Recognition (SER)



[1] Riera, P. et al. (2019). No sample left behind: Towards a comprehensive evaluation of speech emotion recognition system. In *Proc. Workshop on Speech, Music and Mind 2019*.

[2] Chou, H. C. et al. (2022, May). Exploiting annotators' typed description of emotion perception to maximize utilization of ratings for speech emotion recognition. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (pp. 7717-7721). IEEE.

Emotion Co-occurrence Weight Penalty Loss function



(a) Annotations

(b) Co-occurrence frequency matrix

(c) Co-occurrence weight matrix

(d) Penalization matrix

This loss function [1] can **penalize more infrequent emotions that do not co-occur**

We denote the penalization matrix as P

$$\begin{aligned} \mathcal{P}\mathcal{L} &= \sum_{i=1}^N (\mathcal{L}_i \cdot P) \\ &= \sum_{i=1}^N \left(\sum_{j=1}^K \sum_{z=1}^K P_{jz} \cdot f_{loss}(Y_{ij}^T, Y_{ij}^P) \right) \end{aligned}$$

f_{loss} = Binary Cross-Entropy (BCE) in the work

[1] Chou, H. C., Lee, C. C., & Busso, C. (2022). Exploiting co-occurrence frequency of emotions in perceptual evaluations to train a speech emotion classifier. In *Proc. Interspeech* (Vol. 2022).

Class-balanced Sigmoid Cross-entropy Loss:

$$\mathcal{L}_{CBL} = \sum_{j=1}^K \left(\frac{1-\beta}{1-\beta^{n_j}} \cdot \mathcal{L}_{BCE}^{(j)} \right),$$

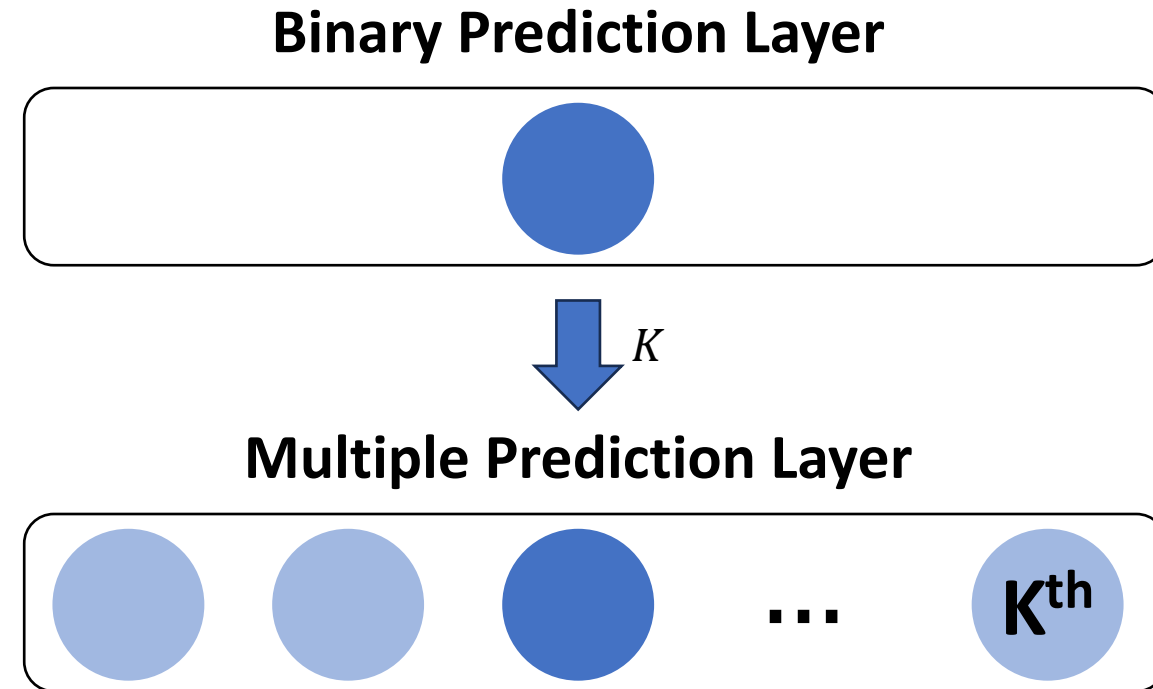
where

- K: number of emotional classes;
- $\beta \in (0, 1]$ hyperparameter ;
- n_j : number of positive samples in j^{th} emotion classes in the “train set”

Add a weighting factor to adjust the values of the used loss function based on the inverses of the class frequency

[1] Cui, et al. (2019). Class-balanced loss based on effective number of samples. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 9268-9277).

Let's consider multi-label tasks as multiple (K) binary classification tasks

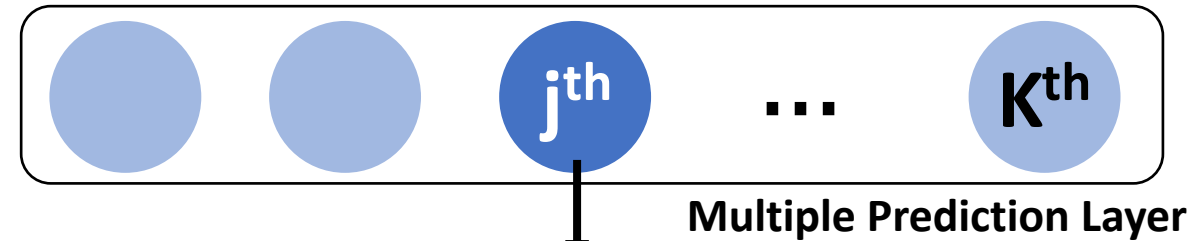


Multi-Label Temperature Scaling Calibration (1/2)

Step 1: Access probabilities of samples in the development set (N is numbers of samples)

- Extract the prediction probabilities from the pre-trained models by freezing the models' weights
- Converting vectors for the emotion j into $N \times 2$, denoted as $Z^P(j)$
- Confidence value (c) of each prediction is the maximum probability of the predictions

Freeze Whole Model's Weights



$$Z_0^P(j)$$

$$Z_1^P(j) = 1 - Z_0^P(j)$$

$$Z^P(j) = [Z_0^P(j) \ Z_1^P(j)]_{N \times 2}$$

$$c(j) = \max(Z_0^P(j), Z_1^P(j))$$

Multi-Label Temperature Scaling Calibration (2/2)

Step 2: Get the logits matrix by using the LogSoftmax activation function

Step 3: Divide the logits vector by a learnable single scalar temperature (T)

- The model learns the optimal value for T by minimizing the negative log-likelihood (NLL) loss on the development set

Step 4: Calculate the calibrated confidences:

$$T(j) \begin{cases} = 1; \text{ maintain the original confidences} \\ > 1; \text{ "smooth" the confidences} \end{cases}$$

$$\begin{aligned} a(j) &= \text{LogSoftmax}(Z^P) \\ &= \log \left(\frac{\exp(Z^P)}{\sum_q \exp(Z_q^P)} \right) \end{aligned}$$

$$\frac{a(j)}{T(j)} \longrightarrow \text{Learnable } T$$

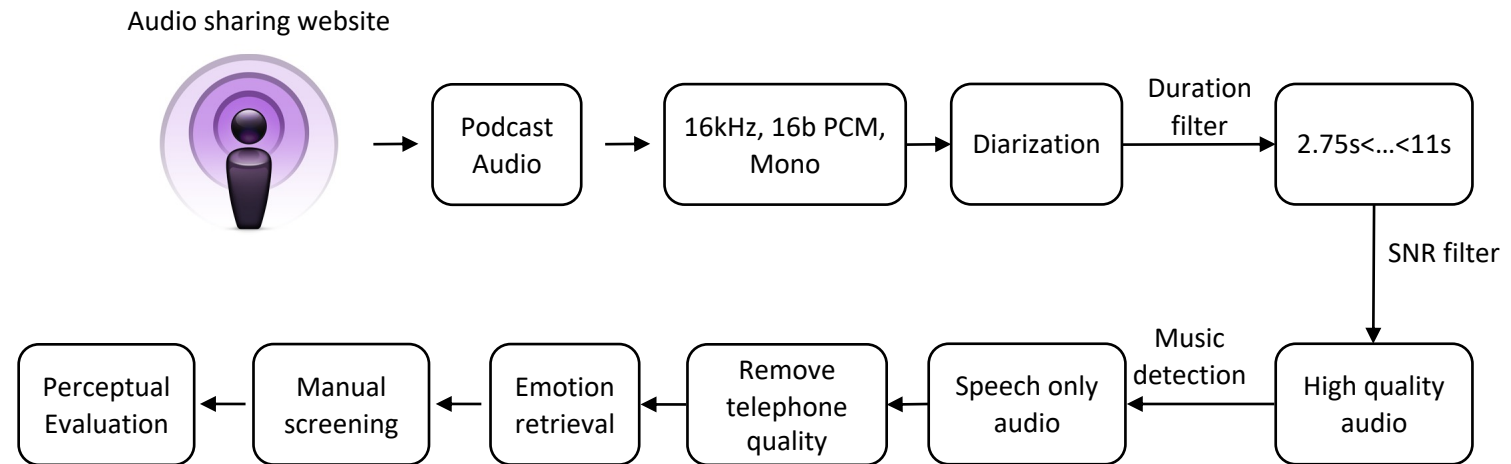
$$\begin{aligned} c(j) &= \text{Softmax} \left(\frac{a(j)}{T(j)} \right) \\ &= \frac{\exp(a(j)/T(j))}{\sum_q \exp(a_q(j)/T(j))} \end{aligned}$$

Audio sentences:

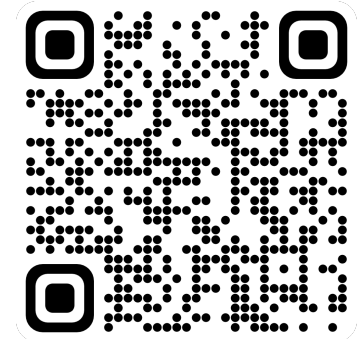
- Train set: 63,076
- Validation set: 10,999
- Test set: 16,903

Emotional Annotations:

- Crowdsource-based protocol
- Every sentence has **more than 5** annotators
- Primary emotion (P) (Single-choice):
 - anger, sadness, happiness, surprise, fear, disgust, contempt, neutral, and ~~other (excluded)~~



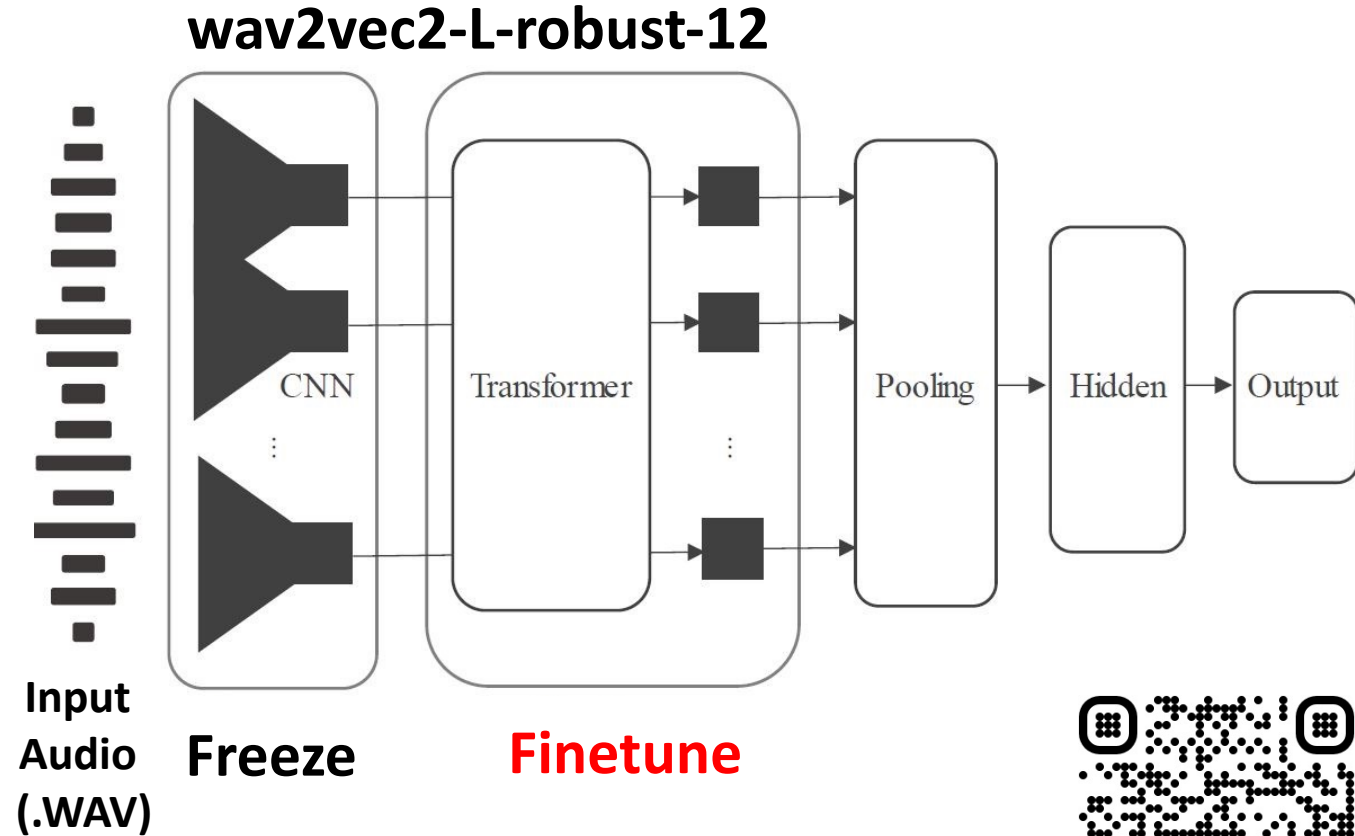
Access the MSP-Podcast



Experiment Setup (SER Framework)

We use the “wav2vec2-L-robust-12” model in [1] and follow the finetuning process:

- Freeze the weights of the *convolutional neural network* (CNN) layers
- Finetune the weights of the 12 transformer layers on the emotion corpus
- Use Adam optimizer with a learning rate of 0.0001 and a batch size of 32



Code of the Original paper[1]



Model evaluation of different systems on primary emotion recognition

1. \mathcal{L}_{BCE} : Binary cross-entropy loss (as Baseline)
2. \mathcal{L}_{CBL} : Consider class-balanced loss
3. \mathcal{L}_P : $(1-\alpha) \cdot \mathcal{L}_{BCE} + \alpha \cdot P\mathcal{L}$
4. $\mathcal{L}_{P+CB} = (1-\alpha) \cdot \mathcal{L}_{CBL} + \alpha \cdot \mathcal{L}_{P+CB}$
5. Apply the multilabel temperature scaling calibration method
 - Can be used for the above No. 1, No.2, No.3, and No. 4 systems

$P\mathcal{L}$ = Emotion co-occurrence weight penalty loss

$\alpha = 0.0, 0.2, 0.5, \text{ or } 0.8$ (we did not optimize α in this work)

[1] Chou, H. C., Lee, C. C., & Busso, C. (2022). Exploiting co-occurrence frequency of emotions in perceptual evaluations to train a speech emotion classifier. In *Proc. Interspeech* (Vol. 2022).

Calibration: Expected Calibration Error (ECE) [1]:

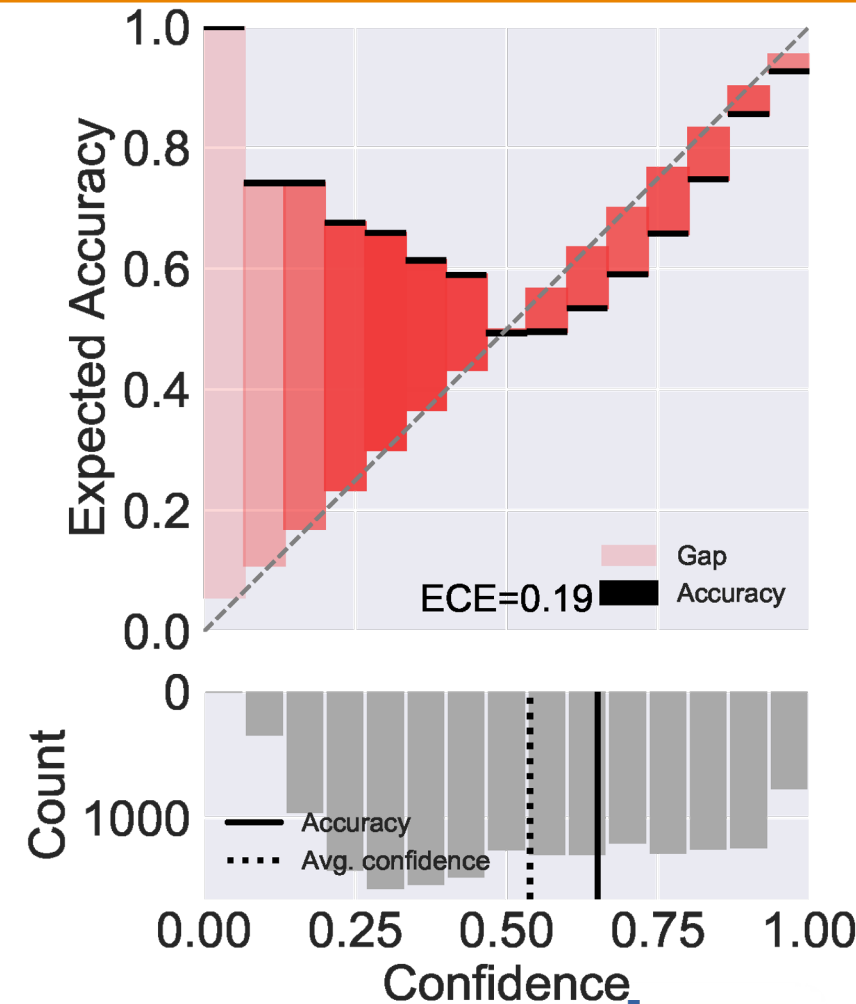
$$ECE = \sum_{b=1}^B \frac{N^b}{N} |accuracy^b - confidence^b|,$$

where

- B = # of bins ($B = 15$, follow [1])
- N^b = # of samples in the b^{th} bin
 - We calculate the ECE scores for 8 emotions and show the averaged ECE as the final scores

Classification Performance:

- Macro-F1 (maF1) - Other metrics in the paper



Average of results:

- Split the original test set into 40 small subsets
- Report the average results in evaluation metrics

Tests of Significance:

- Perform a **two-tailed test**
- Evaluate the statistical significance of all results between the proposed method and baselines
 - Use the symbol * to indicate that the results of the models are statistically significant over the baseline

Research Question (1)

Can CBL and PL improve the performance and calibration of an SER system?

- \mathcal{L}_{CBL} can improve both the classification performance and calibration of SER systems
- PL only can improve the classification performance based on \mathcal{L}_P results
- Use CBL and PL simultaneously can improve the performance and calibration of SER

CBL: class-balanced loss

PL: emotion co-occurrence weight penalty loss

Loss	α	CB	maF1 \uparrow	ECE \downarrow
\mathcal{L}_{BCE}			0.352	0.335
\mathcal{L}_{CBL}		✓	0.367	0.311*
\mathcal{L}_P	0.2		0.320	0.352
	0.5		0.360	0.365
	0.8	-	0.331	0.345
	1.0		0.329	0.351
\mathcal{L}_{P+CB}	0.2		0.401*	0.328
	0.5		0.385*	0.339
	0.8	✓	0.400*	0.329
	1.0		0.371	0.316

Research Question (2)

Can we improve the confidence of the predictions of SER systems without performance drops?

- **CB✓** means the model is calibrated by the proposed multi-label TS calibration method
- CB sustains classification performances
- Improved ECE values with gains between 15.4% and 20%
- **Best performance and calibration achieved with CBL + PL + multi-label TS calibration**

CBL: class-balanced loss

PL: emotion co-occurrence weight penalty loss

Loss	α	CB	maF1 \uparrow	ECE \downarrow	ECE (CB✓) \downarrow	ECE Gain
\mathcal{L}_{BCE}			0.352	0.335	0.276	17.6%
\mathcal{L}_{CBL} :		✓	0.367	0.311*	0.263	15.4%
\mathcal{L}_P	0.2		0.320	0.352	0.292	17.0%
	0.5		0.360	0.365	0.292	20.0%
	0.8	-	0.331	0.345	0.286	17.1%
	1.0		0.329	0.351	0.289	17.7%
\mathcal{L}_{P+CB}	0.2		0.401*	0.328	0.270	17.7%
	0.5		0.385*	0.339	0.277	18.3%
	0.8	✓	0.400*	0.329	0.273	17.0%
	1.0		0.371	0.316	0.266	15.8%

Contribution:

- Multi-label speech emotion classifiers are under-confident
- **Emotion co-occurrence weight penalty function + Class-balanced objective function, + Multi-label calibration** simultaneously can improve performance and calibration of SER models

Results (8-class Primary emotion classification)

- **+2.22%** improvement gain in ECE
- **+13.92%** performance gain in macro- F1 score

Take-Home Message:

- **It is important to calibrate the SER system to improve its confidence and performance**

Thank You



Scan Me for
Paper Full Text

Contact Huang-Cheng Chou:

Email: hc.chou@gapp.nthu.edu.tw

LinkedIn: <https://www.linkedin.com/in/huangchougchou/>



傑出人才發展基金會

Foundation For The Advancement of Outstanding Scholarship



CNS-2016719



NTHU

23