# DISTRIBUTION-BASED EMOTION RECOGNITION IN CONVERSATION

*Wen Wu, Chao Zhang, Philip C. Woodland*

Department of Engineering, University of Cambridge, Trumpington St., Cambridge, UK.

{ww368,cz277,pcw}@eng.cam.ac.uk

## ABSTRACT

Automatic emotion recognition in conversation (ERC) is crucial for emotion-aware conversational artificial intelligence. This paper proposes a distribution-based framework that formulates ERC as a sequence-to-sequence problem for emotion distribution estimation. The inherent ambiguity of emotions and the subjectivity of human perception lead to disagreements in emotion labels, which is handled naturally in our framework from the perspective of uncertainty estimation in emotion distributions. A Bayesian training loss is introduced to improve the uncertainty estimation by conditioning each emotional state on an utterance-specific Dirichlet prior distribution. Experimental results on the IEMOCAP dataset show that ERC outperformed the single-utterance-based system, and the proposed distribution-based ERC methods have not only better classification accuracy, but also show improved uncertainty estimation.

***Index Terms***— automatic emotion recognition, emotion recognition in conversation, Dirichlet prior network, IEMOCAP

## 1. INTRODUCTION

Emotion understanding is a key attribute of conversational artificial intelligence (AI). Although significant progress has been made in developing deep-learning-based automatic emotion recognition (AER) systems over the past several years [1, 2, 3, 4, 5], most studies have focused on modelling and evaluating each utterance separately. However, emotions are known to be dependent on cross-utterance contextual information and persist across multiple utterances in dialogues [6]. This motivates the study of AER in conversation (ERC).

Emotion annotation is challenging due to the fact that emotion is inherently complex and ambiguous, and its expression and perception are highly personal and subjective. This causes uncertainty in the manual references used for emotional data. Although it is common to handle AER as a classification problem based on the majority agreed labels among several annotators [7, 8, 9, 10], it can cause two major problems. First, utterances without majority agreed labels have to be discarded, which makes the dialogue context non-contiguous in both training and test. Second, replacing the (possibly different) original labels from human annotators by the majority vote label also removes the inherent uncertainty associated with emotion perception in the data labelling procedure. To this end, alternative methods to using the majority vote label are required for ERC.

Motivated by these problems, this paper proposes a novel distribution-based framework for ERC, which trains a dialogue-level Transformer model [11] to maximise the probability of generating a

sequence of emotion distributions associated with a sequence of utterances. Each time step of the Transformer represents an utterance in the dialogue, and its corresponding input feature vector is the fusion of audio and text representations for that utterance derived using Wav2Vec 2.0 (W2V2) [12] and bidirectional encoder representations from Transformers (BERT) [13] respectively. The predicted emotion distribution of each utterance relies on all previous predictions as well as the audio and text features in the dialogue. By considering an emotion distribution as a continuous-valued categorical distribution, the original emotion class labels provided by the annotators can be viewed as samples drawn from the underlying true emotion distribution of the utterance. The proposed distribution-based ERC system then learns the true emotion distribution sequence in a conversation given the observed label samples. A novel training loss based on utterance-specific Dirichlet priors predicted by a Dirichlet prior network (DPN) is applied [14], which improves distribution modelling performance by retaining the uncertainty in the original labels. Furthermore, by considering emotion as distributions, no utterances need to be discarded in either training or test, which keeps the dialogue context contiguous.

The rest of the paper is organised as follows. Section 2 provides background to ERC and the modelling of emotion ambiguity. Section 3 introduces distribution-based ERC. Section 4 describes the use of representations derived by self-supervised learning (SSL) and the fusion of audio and text representations. The results and analysis are given in Section 5, followed by conclusions.

## 2. RELATED WORK

### 2.1. Emotion recognition in conversation

Emotion states can be described by discrete emotion categories (*i.e.*, anger, happiness, neutral, sadness, *etc.*) [15] or continuous emotion attributes (*i.e.*, valence-arousal) [16, 17]. This work focuses on classification-based AER using discrete emotion categories. Much work has been published in classification-based AER using deep-learning-based methods [2, 4, 5, 18]. While good recognition performance has been achieved, the focus is on modelling information of each target utterance independently without considering the cross-utterance contextual information. Incorporating such information from both speakers in a dyadic conversation has been shown to improve AER performance [19].

The conversational memory network (CMN) [20] was one of the first ERC approaches that used separate memory networks for both interlocutors participating in a dyadic conversation. Built on the CMN, the interaction-aware attention network (IANN) [21] integrates distinct memories of each speaker using an attention mechanism. Recently, the graph convolutional network was introduced to explore the relations between utterances in a dialogue [22, 23, 24] where the representation of each utterance is treated as nodes and the

relations between the utterances are the edges. Most of these models integrated a fixed-length context rather than the complete dialogue history. Moreover, while contextual information was incorporated by including features of context utterances as inputs, the dependency of the output emotion states was not considered.

Several alternative structures have been proposed in response to the above two issues. The DialogRNN [25] approach uses a hierarchical recurrent neural network framework to recurrently model the emotion of the current utterance by considering the speaker states and the emotions of preceding utterances. The Transformer encoder structure has been adopted for ERC with a data augmentation method based on random utterance concatenation [26]. Another approach is to introduce an extra dialogue-level model on top of the single-utterance-based AER classifier. In the emotion interaction and transition method [27], the emotion probability of each utterance is re-estimated using the previous utterance and currently estimated posteriors using an additional long short-term memory network. The dialogical emotion decoding (DED) method [28] treats a dialogue as a sequence and consecutively decodes the emotion states of each utterance over time with a given recognition engine.
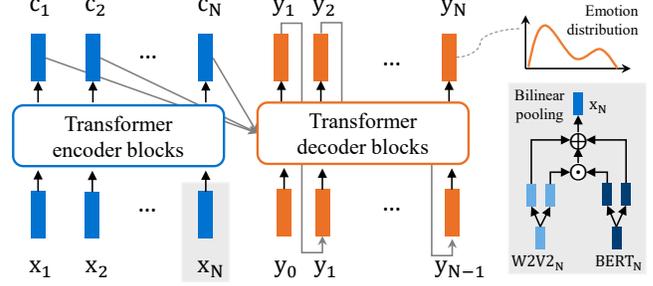
## 2.2. Modelling emotion ambiguity

The methods reviewed in Section 2.1 only used utterances that have majority agreed emotion labels. However, emotion is inherently ambiguous and its perception is highly subjective, which leads to a large degree of uncertainty in labels. Different human annotators may assign different emotion labels to the same utterance and a considerable amount of data does not have majority agreed labels from the human annotators. Majority voting results among the original labels are usually taken as the ground-truth label by AER datasets [7, 8, 9, 10]. Data without majority agreed labels are usually excluded from both training and test in classification-based AER, as they cannot be evaluated without ground-truth labels. In ERC, data exclusion can cause non-contiguous dialogue contexts for both training and test. More importantly, the majority voting strategy considerably changes the uncertainty properties of the emotion states of the speaker [14]. Therefore, alternative methods for emotion state modelling are required. Soft labels have been used as targets in single-utterance-based AER [29, 30], which averages the original emotion class labels provided by the annotators. Such soft labels can be interpreted as the intensities of each emotion class and can allow all utterances to have a training label. Despite the use of soft labels, the systems were evaluated based on classification accuracy with the utterances with majority agreed labels, which results in a mismatch between the training loss and the evaluation metric [31]. In this work, we propose a distribution-based method that allows the use of all utterances and all original labels for ERC.

## 3. DISTRIBUTION-BASED ERC

### 3.1. ERC as a sequence-to-sequence problem

Consider a dialogue with $N$ utterances, let $\mathbf{x}_n$ and $\mathbf{y}_n$ be an utterance and its emotion state in terms of a probability distribution, ERC can be formulated as a special sequence-to-sequence problem, in which the input utterance sequence $\mathbf{x}_{1:N}$ and output emotion state sequence $\mathbf{y}_{1:N}$ have equal lengths and each output distribution represents the emotion state of the corresponding input utterance. Training can be performed by maximising $p(\mathbf{y}_{1:N}|\mathbf{x}_{1:N})$, the likelihood of generating the emotion states based on the input utterances. Based on the



**Fig. 1**. Schematic of the proposed distribution-based ERC system with a Transformer, where $\mathbf{x}_n$ is an utterance and $\mathbf{y}_n$ is the corresponding emotion distribution. Bilinear pooling is used to fuse the audio and text features derived from W2V2 and BERT.

chain rule, $p(\mathbf{y}_{1:N}|\mathbf{x}_{1:N})$ can be calculated efficiently as:

$$p(\mathbf{y}_{1:N}|\mathbf{x}_{1:N}) = p(\mathbf{y}_1|\mathbf{x}_{1:N}) \prod_{n=1}^{N} p(\mathbf{y}_{n+1}|\mathbf{y}_{1:n}, \mathbf{x}_{1:N}). \quad (1)$$

Eqn. (1) differs from single-utterance-based AER [1, 2, 3, 4, 5] by conditioning $\mathbf{y}_{n+1}$ not only on $\mathbf{x}_{1:N}$ but also on $\mathbf{y}_{1:n}$, which reflects the fact that emotional states often persist across multiple utterances in a dialogue [6].

### 3.1.1. A Transformer for online ERC

In Eqn. (1), $p(\mathbf{y}_{n+1}|\mathbf{y}_{1:n}, \mathbf{x}_{1:N})$ not only depends on the current and previous utterances $\mathbf{x}_{1:n+1}$, but also on future utterances $\mathbf{x}_{n+2:N}$, which is not suitable for online AER applications. Hence, an independence approximation between $\mathbf{y}_{n+1}$ and $\mathbf{x}_{n+2:N}$ can be made:

$$p(\mathbf{y}_{1:N}|\mathbf{x}_{1:N}) \approx p(\mathbf{y}_1|\mathbf{x}_1) \prod_{n=1}^{N} p(\mathbf{y}_{n+1}|\mathbf{y}_{1:n}, \mathbf{x}_{1:n+1}), \quad (2)$$

which is used throughout this paper.

In this paper, Eqn. (2) is implemented with a Transformer encoder-decoder model [11], the schematic of the proposed system is shown in Fig. 1. Teacher-forcing [32] is commonly used when training an auto-regressive decoder structure where the output of the current time step depends on outputs of previous time steps. When making a prediction $\mathbf{y}_n$ at a time step $n$, teacher-forcing uses the ground-truth label history $\mathbf{t}_{1:n-1}$ during training, and uses the previous predictions output by the model $\hat{\mathbf{y}}_{1:n-1}$ during test. Teacher-forcing forces the decoder to over-fit to the ground-truth-label-based history and leads to a discrepancy between training and test. Such a discrepancy can yield errors that propagate and accumulate quickly along the generated sequence.

### 3.1.2. Avoid over-fitting to oracle history with scheduled sampling

To alleviate the discrepancy caused by teacher-forcing, scheduled sampling [33] is used during training in this paper. A teacher-forcing ratio $\epsilon_i$ is introduced to randomly decide, during training, whether to use $\mathbf{t}_{n-1}$ or $\hat{\mathbf{y}}_{n-1}$:

$$\mathbf{y}_{n-1} = \begin{cases} \mathbf{t}_{n-1}, & p_{\text{tf}} \leq \epsilon_i \\ \hat{\mathbf{y}}_{n-1}, & p_{\text{tf}} > \epsilon_i \end{cases} \quad (3)$$

where $p_{\text{tf}}$ is sampled from a uniform distribution between 0 and 1 ($p_{\text{tf}} \sim \mathbf{U}(0, 1)$), $i$ denotes the $i^{th}$ mini-batch. $\epsilon_i$ gradually decreases

as the training progresses, which changes the training process from a fully guided scheme based on previous ground-truth labels towards a less guided scheme based on previous model outputs. Commonly used schedules decrease $\epsilon_i$ as a function of $i$, which include a linear decay, an exponential decay and a inverse sigmoid decay [33].

## 3.2. Emotion distribution modelling using DPN

Denote the underlying true emotion distribution of an utterance $\mathbf{x}$ as a categorical distribution $\boldsymbol{\mu} = [p(\omega_1|\boldsymbol{\mu}), \ldots, p(\omega_K|\boldsymbol{\mu})]^{\mathrm{T}}$, where $K$ is the number of emotion classes. Emotion class labels $\omega_k$ from human annotators are samples drawn from this categorical distribution. Soft labels, as an approximation to the underlying true distribution, correspond to the maximum likelihood estimate (MLE) of $\boldsymbol{\mu}$ given the observed label samples $\{\boldsymbol{\mu}^{(1)}, \ldots, \boldsymbol{\mu}^{(M)}\}$:

$$\bar{\boldsymbol{\mu}} = \arg\max_{\boldsymbol{\mu}} \ln p(\boldsymbol{\mu}^{(1)}, \ldots, \boldsymbol{\mu}^{(M)}|\boldsymbol{\mu}) = \frac{1}{M}\sum_{m=1}^{M}\boldsymbol{\mu}^{(m)}. \quad (4)$$

Although soft labels enable contiguous dialogue contexts to be modelled by ERC, the MLE is a good approximation to the true distribution only if a large number of original labels are available for each utterance, which cannot usually be satisfied in practice due to labelling difficulty and cost. Here, we use the Dirichlet prior network (DPN) [34, 35, 14] to resolve the label sparsity issue, which is a Bayesian approach modelling $p(\boldsymbol{\mu}|\boldsymbol{x})$ by predicting the parameters of its Dirichlet prior distribution.

### 3.2.1. Emotion recognition with a Dirichlet prior

The Dirichlet distribution, as the conjugate prior of the categorical distribution, is parameterised by its concentration parameter $\boldsymbol{\alpha} = [\alpha_1, \ldots, \alpha_K]^{\mathrm{T}}$. A Dirichlet distribution $\mathrm{Dir}(\boldsymbol{\mu}|\boldsymbol{\alpha})$ is

$$\mathrm{Dir}(\boldsymbol{\mu}|\boldsymbol{\alpha}) = \frac{\Gamma(\alpha_0)}{\prod_{k=1}^{K}\Gamma(\alpha_k)}\prod_{k=1}^{K}\mu_k^{\alpha_k-1}, \quad (5)$$

$$\alpha_0 = \sum_{k=1}^{K}\alpha_k, \;\; \alpha_k > 0,$$

where $\Gamma(\cdot)$ is the gamma function defined as

$$\Gamma(\alpha_k) = \int_0^{\infty} z^{\alpha_k-1}e^{-z}\,dz. \quad (6)$$

In the Dirichlet process, given $\boldsymbol{\alpha}$, a categorical emotion distribution $\boldsymbol{\mu}$ is drawn from $\mathrm{Dir}(\boldsymbol{\mu}|\boldsymbol{\alpha})$, and an emotion class label $\omega_k$ is sampled from $\boldsymbol{\mu}$.

### 3.2.2. DPN for emotion recognition

A DPN is a neural network modelling $p(\boldsymbol{\mu}|\mathbf{x}, \boldsymbol{\Lambda})$ by predicting the concentration parameter $\boldsymbol{\alpha}$ of the Dirichlet prior, where $\boldsymbol{\Lambda}$ is the collection of model parameters. For each utterance $\mathbf{x}$, the DPN predicts $\boldsymbol{\alpha} = \exp[f_{\boldsymbol{\Lambda}}(\boldsymbol{x})]$, where $f_{\boldsymbol{\Lambda}}(.)$ is the DPN model. Note that the predicted $\boldsymbol{\alpha}$ depends on, and is specific to, each input utterance $\mathbf{x}$. By predicting $\boldsymbol{\alpha}$ separately for each utterance, the DPN makes the Dirichlet prior "utterance-specific" and thus suitable for ERC tasks.

The predictive distribution of the DPN is the expected categorical distribution under $\mathrm{Dir}(\boldsymbol{\mu}|\boldsymbol{\alpha})$ [34]:

$$p(\omega_k|\mathbf{x}, \boldsymbol{\Lambda}) = \mathbb{E}_{p(\boldsymbol{\mu}|\mathbf{x}, \boldsymbol{\Lambda})}[p(\omega_k|\boldsymbol{\mu})]$$
$$= \frac{\alpha_k}{\sum_{k'=1}^{K}\alpha_{k'}} = \mathrm{softmax}[f_{\boldsymbol{\Lambda}}(\boldsymbol{x})]_k, \quad (7)$$

which makes the DPN a normal neural network model with a softmax output activation function during test.

DPN training is performed by maximising the likelihood of sampling the original labels (one-hot categorical distributions) from their relevant utterance-specific Dirichlet priors. Given an utterance $\mathbf{x}$ with $M$ original labels $\{\boldsymbol{\mu}^{(1)}, \ldots, \boldsymbol{\mu}^{(M)}\}$, a DPN is trained to minimise the negative log likelihood

$$\mathcal{L}_{\mathrm{dpn}} = -\frac{1}{M}\sum_{m=1}^{M}\ln p(\boldsymbol{\mu}^{(m)}|\mathbf{x}, \boldsymbol{\Lambda}) \quad (8)$$
$$= -\frac{1}{M}\sum_{m=1}^{M}\ln\mathrm{Dir}(\boldsymbol{\mu}^{(m)}|f_{\boldsymbol{\Lambda}}(\mathbf{x})).$$

In contrast to soft labels that retain only the proportion of occurrences of each emotion class, the DPN preserves the information about each single occurrence of the emotion classes and also resolves the label sparsity issue with Dirichlet priors.

However, $\mathcal{L}_{\mathrm{dpn}}$ is not directly applicable to an auto-regressive ERC system. This is because, in such a system, the targets of the previous time step are required for training (whether using teacher-forcing or scheduled sampling). In a DPN system, the output of the network is the hyperparameter $\alpha$ and the targets associated with $\alpha$ of the previous time step is unknown. Therefore, $\mathcal{L}_{\mathrm{dpn}}$ was added as an extra term to the Kullback-Leibler (KL) divergence $\mathcal{L}_{\mathrm{kl}}$ between the soft labels and the DPN predictive distributions. That is,

$$\mathcal{L}_{\mathrm{kl}} = \sum_{k=1}^{K} -\bar{\mu}_k \ln p(\omega_k|\boldsymbol{x}, \boldsymbol{\Lambda}) + \sum_{k=1}^{K}\bar{\mu}_k\ln\bar{\mu}_k \quad (9)$$

$$\mathcal{L}_{\mathrm{dpn\text{-}kl}} = \mathcal{L}_{\mathrm{dpn}} + \lambda\,\mathcal{L}_{\mathrm{kl}}. \quad (10)$$

The targets of the previous time step are the soft labels. This is a major difference from [14], which treats the DPN loss as the main loss to model the uncertainty of emotions for each utterance independently without taking the previous emotion predictions into account.

## 3.3. Evaluating distribution-based ERC

Since classification accuracy cannot be applied to the utterances without majority agreed labels, it is no longer suitable to be the primary measure for evaluating distribution-based ERC system. The area under the precision-recall curve (AUPR) is used as an alternative metric [14] at test-time. A precision-recall (PR) curve is obtained by calculating the precision and recall for different decision thresholds where the $x$-axis of a PR curve is the recall, the $y$-axis is the precision. The AUPR is the average of precision across all recall values computed as the area under the PR curve. Compared to classification accuracy, AUPR can be applied to all test utterances and also quantify the model's ability to estimate uncertainty.

In this paper, the curve is drawn by detecting utterances without majority agreed labels based on the model prediction. Utterances with majority agreed labels are selected as positive class and utterances without majority agreed labels are chosen as the negative class. Two threshold measures representing the confidence encapsulated in the prediction can be used as the threshold for AUPR:

- The entropy of the predictive distribution (Ent.), defined as $-\sum_{k=1}^{K} p(\omega_k|\mathbf{x}, \boldsymbol{\Lambda})\ln p(\omega_k|\mathbf{x}, \boldsymbol{\Lambda})$ that measures the flatness of the emotion distribution, where $\mathbf{x}$ is an utterance, $\omega_k$ is the $k$-th class and $\boldsymbol{\Lambda}$ is the model.

- The max probability (Max.P), $\max_k p(\omega_k|\mathbf{x}, \boldsymbol{\Lambda})$ measuring the confidence of the predicted emotion class.

## 4. SSL REPRESENTATIONS FOR AUDIO AND TEXT

Representations extracted from pre-trained universal models have recently attracted wide attention. These models are trained on a wide range of data at scale and can be fine-tuned to various downstream tasks and are sometimes referred to as foundation models [36]. Self-supervised learning (SSL) is one of the most common approaches to train a foundation model as it does not require any labels from human annotators. It uses information extracted only from the input data itself in order to learn representations useful for downstream tasks, thus allowing the use of a large amount of unlabelled data for model training.

Models pre-trained by SSL have achieved great successes in natural language processing (*e.g.* BERT [13], GPT-2 [37] and GPT-3 [38]) and computer vision (*e.g.* ViT [39] and iGPT [40]), and have attracted increasing attention in speech processing [41] (*e.g.* ACPC [42], W2V2 [12], Hubert [43], and WavLM [44] *etc.*). The large amount of unlabelled data leveraged by SSL can cover many linguistic and para-linguistic phenomena, and therefore could help to alleviate the data sparsity issue in AER [45, 46, 47, 48]. This paper used two SSL models: W2V2 [12] for the audio modality and BERT [13] for the text modality[1].

### 4.1. Features derived from W2V2 and BERT

W2V2 is a type of contrastive [49] SSL model which learns representations by distinguishing a target sample (positive) from distractor samples (negatives) given an anchor representation. It takes as input a waveform and uses a convolutional feature encoder followed by a transformer network. This paper uses the "wav2vec2-base" model[2] which contains 12 transformer blocks with model dimension 768, inner dimension 3,072 and 8 attention heads and is pre-trained using 960 hours of audio from Librispeech corpus [50]. W2V2 representations were extracted from the output of the last Transformer block and averaged across each utterance.

BERT is a type of predictive SSL model which learns representations by predicting the masked tokens in a sentence. This paper uses the "bert-base-uncased" model[3] which contains 12 transformer blocks with a model dimension of 768, inner dimension 3,072 and 12 attention heads and was pre-trained using a large amount of text.

### 4.2. Bilinear-pooling-based feature fusion

Bilinear pooling [51] is a commonly used approach for the expressive fusion of multimodal representations [52], which models the multiplicative interactions between all possible element pairs. It computes the outer product of the $Q$-dimensional (-dim) audio and text representations, $\mathbf{e}_1$ and $\mathbf{e}_2$, into a $Q^2$-dim joint representation and then projects it into an $O$-dim space with a linear transform. In practice, bilinear pooling often suffers from a data sparsity issue caused by the high dimensionality of $Q^2$. Therefore, decomposition techniques are usually required in order to estimate the associated parameters properly and efficiently. In this paper, the W2V2 and BERT features were fused using a modified multimodal low-rank bilinear attention network with shortcut connections [53]:
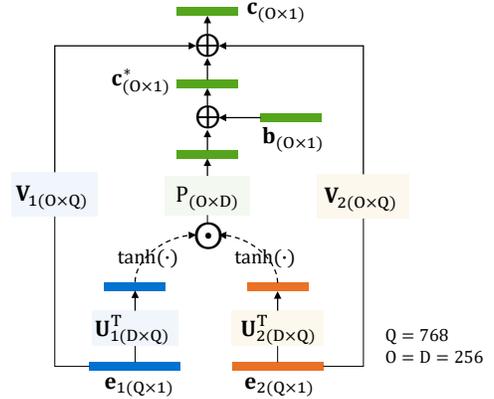
$$\mathbf{c}^* = \mathbf{P}(\tanh(\mathbf{U}_1^{\mathrm{T}}\mathbf{e}_1) \odot \tanh(\mathbf{U}_2^{\mathrm{T}}\mathbf{e}_2)) + \mathbf{b}$$

$$\mathbf{c} = \mathbf{c}^* + \mathbf{V}_1\mathbf{e}_1 + \mathbf{V}_2\mathbf{e}_2.$$

---

[1] The proposed system can be viewed as a dialogue-level audio-text multimodal adaptor for emotions in the foundation model paradigm.

[2] https://huggingface.co/facebook/wav2vec2-base

[3] https://huggingface.co/bert-base-uncased



**Fig. 2**. Schematic of bilinear pooling with shortcut combination. $\odot$ represents the Hadamard product and $\oplus$ represents the element-wise addition of vectors.

The process is illustrated in Fig. 2. In this paper, $\mathbf{e}_1$, $\mathbf{e}_2$ are the W2V2 and BERT derived vectors to be combined, and are both 768-dim. $D$ and $O$ are both 256-dim, $\mathbf{c}$ is the 256-dim combined vector, $\mathbf{U}_1$, $\mathbf{U}_2$ are both 768×256-dim matrices, $\mathbf{b}$ is a 256-dim bias vector, $\mathbf{P}$ is a 256×256-dim linear projection, $\mathbf{V}_1$, $\mathbf{V}_2$ are both 256×768-dim, and $\odot$ is the Hadamard product.

## 5. EXPERIMENTS

### 5.1. Experimental setup

#### 5.1.1. Dataset

The IEMOCAP [7] corpus is used in this paper, which is one of the most widely used datasets for verbal emotion classification. It consists of 5 dyadic conversational sessions performed by 10 professional actors with a session being a conversation between two speakers. There are in total 151 dialogues including 10,039 utterances. Each utterance was annotated by three human annotators for emotion class labels (neutral, happy, sad, and angry *etc.*). Each annotator was allowed to tag more than one emotion category for each sentence if they perceived a mixture of emotions, giving utterances 3.12 labels on average. Ground-truth labels were determined by majority vote Following prior work [1, 2, 3], the reference transcriptions from IEMOCAP were used for the text modality. Leave-one-session-out 5-fold cross validation (5-CV) was performed and the average results are reported.

#### 5.1.2. Data augmentation

In the dialogue-level ERC system, the number of training samples is equal to the number of dialogues in the dataset, which is often very limited (*e.g.* 151 in IEMOCAP). To mitigate training data sparsity issues, sub-sequence randomisation [54] was used to augment the training data, which samples sub-sequences $(\mathbf{x}_{s:e}, \mathbf{y}_{s:e})$ from each full training sequence as the augmented training samples, where $s$ and $e$ are the randomly selected start and end utterance indexes.

#### 5.1.3. Training specifications

The Transformer architecture [11] used for ERC contains 4 encoder blocks and 4 decoder blocks with a dimension of 256. The multi-head attention contains 4 heads. Masking was applied to ensure

**Table 1**. 5-fold CV classification results (mean± standard deviation across folds) for 4-way utterance and dialogue baseline systems on IEMOCAP. IANN [21] and DED [28] are ERC methods using audio features only without pre-trained encoders.

| System | %WA | %UA |
|---|---|---|
| utterance-W2V2ft | 68.71±2.60 | 69.99±3.91 |
| dialogue-W2V2ft | 70.18±4.83 | 71.32±4.06 |
| dialogue-BERT | 68.57±3.50 | 67.56±2.94 |
| dialogue-W2V2ft+BERT | 74.87±3.77 | 74.59±3.50 |
| IANN [21] | 64.7 | 66.3 |
| DED [28] | 69.0 | 70.1 |

that predictions for the current utterance depend only on previous input utterances and known outputs. Sinusoidal positional embeddings [11] were added to the input features. A dropout rate of 10% was applied to all parameters. The model was implemented using PyTorch. Scheduled sampling with an exponential scheduler was used during training. The Adam optimiser was used with a variable learning rate with linear warm-up for the first 2,000 training updates and then linearly decreasing.[4]
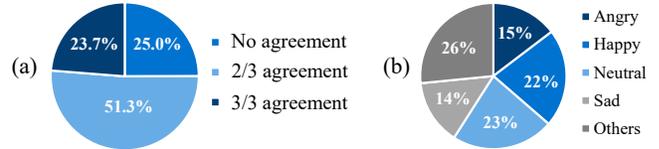
### 5.2. AER with 4-way classification

The most common setup on IEMOCAP [2, 4, 18, 21, 28, 22] only uses utterances with majority agreed labels belonging to "angry", "happy", "excited" (merged with "happy"), "sad", and "neutral" for a 4-way classification, which results in 5,531 utterances. For comparison, 4-way emotion classification systems using this setup were first built as baselines. Since the test sets are slightly imbalanced between different emotion categories, both weighted accuracy (WA) and unweighted accuracy (UA) are reported for classification experiments. WA corresponds to the overall accuracy while UA corresponds to the mean of class-wise accuracy.

The "wav2vec2-base" model was fine-tuned for single-utterance-based 4-way classification by adding a 128-d fully connected layer and an output layer with softmax activation on top of the pre-trained model. The fine-tuning experiment results are shown as "utterance-W2V2ft" in Table 1. The Transformer ERC model was then trained using the fine-tuned W2V2 features ("dialogue-W2V2ft"), BERT features ("dialogue-BERT"), and the fusion of BERT and the fine-tuned W2V2 features ("dialogue-W2V2ft+BERT") as input. Comparing "utterance-W2V2ft" and "dialogue-W2V2ft" in Table 1, dialogue-based AER performs better than single-utterance-based AER. The fusion of audio and text features further improves the performance.

### 5.3. IEMOCAP data analysis and motivation for distribution-based systems

Statistics for the IEMOCAP corpus are summarized in Fig. 3. The 4-way classification setup discards 45% of the data in IEMOCAP that belongs to the following two categories:

- Utterances without majority agreed emotion labels (2,507 utterances);
- Utterances whose majority agreed labels do not belong to the selected four emotion classes in the 4-way setup (2,001 utterances).



**Fig. 3**. Distribution of data in IEMOCAP. (a) Proportion of annotators agreeing on the label. (b) Ground-truth of utterances with unique majority labels.

This strategy not only causes a loss of nearly half of the emotion data, which are highly valuable, but also causes the dialogue contexts modelled by the Transformer to be non-contiguous. Furthermore, among the utterances with majority agreed labels, only 31.6% have all annotators agreed on the same emotion class label. When majority voting is applied, an utterance with labels "happy", "happy", "happy" and an utterance with labels "happy", "happy", "sad" have the same ground-truth label "happy", even though an annotator has assigned a different label to the latter utterance. The use of majority voting therefore changes the true emotion distributions and causes the inherent uncertainty that exists in the annotations to be discarded.

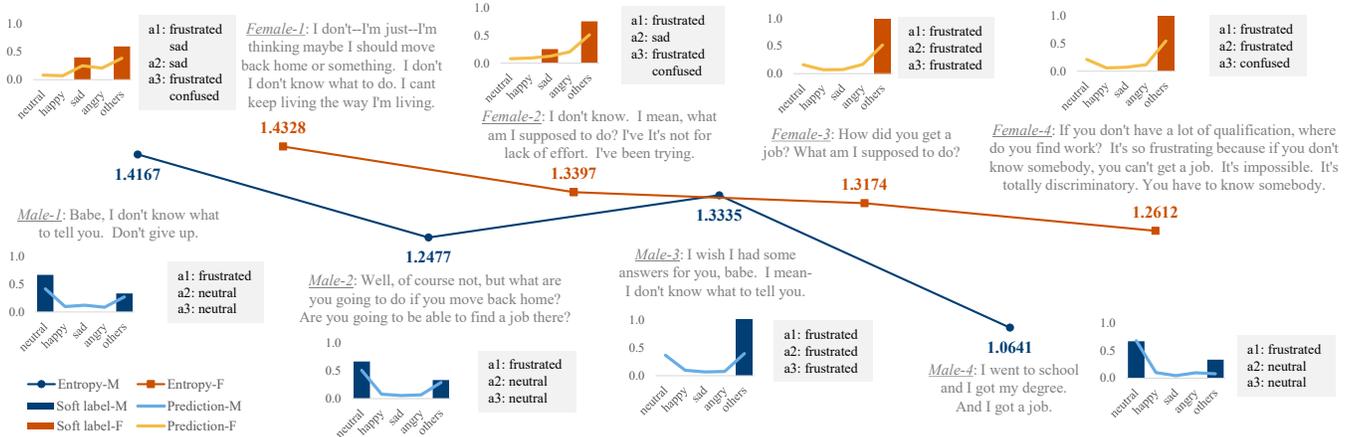### 5.4. ERC with Distribution modelling

In this section, the training targets for each utterance were revised from a single discrete label to a continuous-valued categorical distribution over five emotion classes. The five classes correspond to the previous four emotion classes plus an extra class "others" that includes all the other emotions that exist in IEMOCAP. This not only allows all data in IEMOCAP to be used for ERC, regardless of whether majority agreed label exists and whether it belongs to the 4 classes, but also avoids the problem that not all of the original labels are represented by majority voting. Three systems were evaluated:

- HARD: A system trained by 5-way classification. When training the HARD system, all utterances in a dialogue were taken as input while the loss was only computed for utterances that have majority agreed labels.
- SOFT: A system trained by minimising $\mathcal{L}_{kl}$ in Eqn. (9).
- DPN-KL: A system trained by minimising the combined loss $\mathcal{L}_{dpn-kl}$ in Eqn. (10) with $\lambda = 20$.

All systems were first evaluated by 5-way classification accuracy on test utterances with majority agreed labels, as well as by 4-way classification accuracy on test utterances whose majority agreed labels belong to the 4 classes, and then evaluated by the average AUPR (Max.P) and AUPR (Ent.) on all test utterances.

As shown by the results in Table 2, both the SOFT and DPN-KL systems outperform the HARD system on all evaluation metrics. A possible explanation lies in the fact that emotion assignment errors in the hard classification setup (HARD system) are more likely to propagate through the dialogue with the auto-regressive Transformer decoder. The DPN-KL system produces the highest AUPR among all of the systems, which shows that it has the best performance in modelling emotion distributions as it gives the best prediction of utterances without majority agreed labels. For the convenience of visualisation, the 5th fold (trained on Session 1-4 and tested on Session 5) was taken as an example and the PR curves of test utterances for all three systems are shown in Fig. 5. It can be seen that the
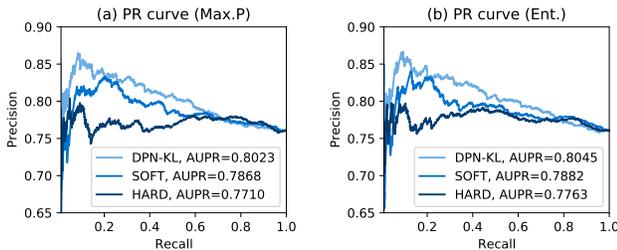
---

[4]Code availble: https://github.com/W-Wu/ERC-SLT22

**Fig. 4**. Entropy of the predicted emotion distribution of each utterance in a sub-dialogue. The DPN-KL system trained on Session 1-4 was used. For each sentence, the bar chart shows the soft label and the line on the bar chart shows the prediction. Labels provided by the three annotators are shown in the grey box, with "a1" referring to the first annotator *etc*. ("frustrated" and "confused" are merged into the 5-th class "others").

**Table 2**. 5-fold CV results for the proposed distribution-based ERC method on IEMOCAP. Highest value in each row shown in bold font.

| Metric | HARD | SOFT | DPN-KL |
|---|---|---|---|
| 5-way %WA | 59.99±4.14 | 62.48±3.70 | **63.63**±2.43 |
| 5-way %UA | 58.83±3.51 | 62.54±3.98 | **63.12**±3.30 |
| 4-way %WA | 73.01±3.79 | 77.46±2.49 | **77.83**±2.07 |
| 4-way %UA | 72.57±3.23 | **78.16**±2.88 | 78.12±2.60 |
| %AUPR (MaxP) | 76.63±2.09 | 78.02±1.32 | **80.72**±1.67 |
| %AUPR (Ent.) | 76.99±2.24 | 78.63±1.27 | **81.17**±2.02 |



**Fig. 5**. PR curves of the three systems using (a) Max.P and (b) Ent. as the uncertainty measures. The tests were performed on Session 5.

DPN-KL system gives the largest area under the PR curve, showing its superior uncertainty estimation performance.

### 5.5. Analysis

This section gives a case study to better understand uncertainty variation of emotion distributions in a dialogue. Fig. 4 shows the trend of uncertainty change in emotion estimation (measured by the entropy of the predicted emotion distribution) in a sub-dialogue selected from the dialogue "Ses05F_impro04" in IEMOCAP Session 5 between a female and a male speaker. The results were produced by the DPN-KL system trained on Session 1-4. For each sentence, the bar chart shows the soft label and the line on the bar chart shows the prediction. Labels provided by the three annotators are shown in the grey box.

From Fig. 4, utterance *Female-1* has two annotators that each provided two labels, indicating the uncertainty of the emotional content of the utterance. The uncertainty of emotion estimation is reflected by the high entropy. Although utterances *Female-2* and *Female-3* have the same emotion class labels found by majority voting, their underlying emotion distributions are different. As the dialogue progressed (from *Female-1* to *Female-3*), the annotators gradually became more certain that the female speaker got frustrated (shown by the soft labels), and the predicted distribution changed accordingly. Given the same label samples (*i.e., Female-3* and *Female-4*; *Male-1*, *Male-2* and *Male-4*), the entropy decreases as the dialogue progressed, indicating that the model is becoming more confident about its predictions. The reductions of uncertainty in this example demonstrates the advantage of using cross-utterance contextual information in our proposed distribution-based ERC framework. Moreover, another example is the emotional shift that occurs from utterance *Male-2* to *Male-3*. Due to emotional inertia, the model predicts a higher probability of "others" while still retaining some probability for "neutral". The larger entropy reveals that the model is uncertain about this prediction.

## 6. CONCLUSION

In this paper, we propose a distribution-based ERC framework, which formulates ERC as a special sequence-to-sequence problem for distribution estimation. The emotion state of an utterance is represented by a categorical distribution which depends on the context information and emotion states of the previous utterances in the dialogue. Each input vector of a dialogue sequence input to the Transformer dialogue model is formed by fusing representations extracted from SSL models W2V2 and BERT, which makes the system also a dialogue-level audio-text multimodal task adaptor for AER. The system is trained by minimising the KL divergence combined with the DPN loss which conditions the categorical distribution on an utterance-specific Dirichlet prior distribution, which is evaluated by AUPR with the task of detecting utterances without majority agreed labels. This approach not only allows utterances without majority agreed labels to be used, but also leads to better performance in modelling the uncertainty variations in ERC.

# 7. REFERENCES

[1] Samarth Tripathi., Sarthak Tripathi, and Homayoon Beigi, "Multi-modal emotion recognition on IEMOCAP dataset using deep learning," *arXiv preprint 1804.05788*, 2018.

[2] N. Majumder, D. Hazarika, A. Gelbukh, E. Cambria, and S. Poria, "Multimodal sentiment analysis using hierarchical fusion with context modeling," *Knowledge-Based Systems*, vol. 161, pp. 124–133, 2018.

[3] Soujanya Poria, Navonil Majumder, Devamanyu Hazarika, Erik Cambria, Alexander Gelbukh, and Amir Hussain, "Multimodal sentiment analysis: Addressing key issues and setting up the baselines," *IEEE Intelligent Systems*, vol. 33, no. 6, pp. 17–25, 2018.

[4] Wenjing Han, Huabin Ruan, Xiaomin Chen, Zhixiang Wang, Haifeng Li, and Björn Schuller, "Towards temporal modelling of categorical speech emotion recognition," in *Proc. Interspeech*, Hyderabad, 2018.

[5] Raghavendra Pappagari, Tianzi Wang, Jesus Villalba, Nanxin Chen, and Najim Dehak, "X-vectors meet emotions: A study on dependencies between emotion and speaker recognition," in *Proc. ICASSP*, Barcelona, 2020.

[6] Jerry Suls, Peter Green, and Stephen Hillis, "Emotional reactivity to everyday problems, affective inertia, and neuroticism," *Personality and Social Psychology Bulletin*, vol. 24, no. 2, pp. 127–136, 1998.

[7] C. Busso, M. Bulut, C.-C. Lee, A. Kazemzadeh, E.M. Provost, S. Kim, J.N. Chang, S. Lee, and S.S. Narayanan, "IEMOCAP: Interactive emotional dyadic motion capture database," *Language Resources and Evaluation*, vol. 42, pp. 335–359, 2008.

[8] Ya Li, Jianhua Tao, Björn Schuller, Shiguang Shan, Dongmei Jiang, and Jia Jia, "MEC 2017: Multimodal emotion recognition challenge," in *Proc. ACII Asia*, Beijing, 2018.

[9] Carlos Busso, Srinivas Parthasarathy, Alec Burmania, Mohammed AbdelWahab, Najmeh Sadoughi, and Emily Mower Provost, "MSP-IMPROV: An acted corpus of dyadic interactions to study emotion perception," *IEEE Trans on Affective Computing*, vol. 8, no. 1, pp. 67–80, 2017.

[10] Soujanya Poria, Devamanyu Hazarika, Navonil Majumder, Gautam Naik, Erik Cambria, and Rada Mihalcea, "MELD: A multimodal multi-party dataset for emotion recognition in conversations," *arXiv preprint arXiv:1810.02508*, 2018.

[11] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin, "Attention is all you need," in *Proc. NeurIPS*, Long Beach, 2017.

[12] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli, "Wav2Vec 2.0: A framework for self-supervised learning of speech representations," in *Proc. NeurIPS*, Virtual, 2020.

[13] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proc. NAACL*, Minneapolis, 2019.

[14] Wen Wu, Chao Zhang, Xixin Wu, and Philip C. Woodland, "Estimating the uncertainty in emotion class labels with utterance-specific dirichlet priors," *arXiv preprint arXiv:2203.04443*, 2022.

[15] Hatice Gunes, Björn Schuller, Maja Pantic, and Roddy Cowie, "Emotion representation, analysis and synthesis in continuous space: A survey," in *Proc. FG*, Santa Barbara, 2011.

[16] Harold Schlosberg, "Three dimensions of emotion.," *Psychological review*, vol. 61, no. 2, pp. 81, 1954.

[17] Mihalis A Nicolaou, Hatice Gunes, and Maja Pantic, "Continuous prediction of spontaneous affect from multiple cues and modalities in valence-arousal space," *IEEE Transactions on Affective Computing*, vol. 2, no. 2, pp. 92–105, 2011.

[18] Seunghyun Yoon, Seokhyun Byun, and Kyomin Jung, "Multimodal speech emotion recognition using audio and text," in *Proc. SLT*, Athens, 2018.

[19] Wen Wu, Chao Zhang, and Philip C. Woodland, "Emotion recognition by fusing time synchronous and time asynchronous representations," in *Proc. ICASSP*, Toronto, 2021.

[20] Devamanyu Hazarika, Soujanya Poria, Amir Zadeh, Erik Cambria, Louis-Philippe Morency, and Roger Zimmermann, "Conversational memory network for emotion recognition in dyadic dialogue videos," in *Proc. ACL*, Melbourne, 2018.

[21] Sung-Lin Yeh, Yun-Shao Lin, and Chi-Chun Lee, "An interaction-aware attention network for speech emotion recognition in spoken dialogs," in *Proc. ICASSP*, Brighton, 2019.

[22] Jiaxing Liu, Yaodong Song, Longbiao Wang, Jianwu Dang, and Ruiguo Yu, "Time-frequency representation learning with graph convolutional network for dialogue-level speech emotion recognition," in *Proc. Interspeech*, Brno, 2021.

[23] Deepanway Ghosal, Navonil Majumder, Soujanya Poria, Niyati Chhaya, and Alexander Gelbukh, "DialogueGCN: A graph convolutional neural network for emotion recognition in conversation," in *Proc. EMNLP*, Hong Kong, 2019.

[24] Dong Zhang, Liangqing Wu, Changlong Sun, Shoushan Li, Qiaoming Zhu, and Guodong Zhou, "Modeling both context- and speaker-sensitive dependence for emotion detection in multi-speaker conversations.," in *Proc. IJCAI*, Macao, 2019.

[25] Navonil Majumder, Soujanya Poria, Devamanyu Hazarika, Rada Mihalcea, Alexander Gelbukh, and Erik Cambria, "DialogueRNN: An attentive RNN for emotion detection in conversations," in *Proc. AAAI*, Honolulu, 2019.

[26] Raghavendra Pappagari, Piotr Żelasko, Jesús Villalba, Laureano Moro-Velazquez, and Najim Dehak, "Beyond isolated utterances: Conversational emotion recognition," in *Proc. ASRU*, Cartagena, 2021.

[27] Ruo Zhang, Atsushi Ando, Satoshi Kobashikawa, and Yushi Aono, "Interaction and transition model for speech emotion recognition in dialogue.," in *Proc. Interspeech*, Stockholm, 2017.

[28] Sung-Lin Yeh, Yun-Shao Lin, and Chi-Chun Lee, "A dialogical emotion decoder for speech emotion recognition in spoken dialog," in *Proc. ICASSP*, Barcelona, 2020.

[29] Atsushi Ando, Satoshi Kobashikawa, Hosana Kamiyama, Ryo Masumura, Yusuke Ijima, and Yushi Aono, "Soft-target training with ambiguous emotional utterances for dnn-based speech emotion classification," in *Proc. ICASSP*, Brighton, 2018.

[30] Haytham M. Fayek, Margaret Lech, and Lawrence Cavedon, "Modeling subjectiveness in emotion recognition with deep neural networks: Ensembles vs soft labels," in *Proc. IJCNN*, Vancouver, 2016.

[31] Emily Mower, Angeliki Metallinou, Chi chun Lee, Abe Kazemzadeh, Carlos Busso, Sungbok Lee, and Shrikanth Narayanan, "Interpreting ambiguous emotional expressions," in *Proc. ACII*, Amsterdam, 2009.

[32] Ronald J Williams and David Zipser, "A learning algorithm for continually running fully recurrent neural networks," *Neural computation*, vol. 1, no. 2, pp. 270–280, 1989.

[33] Samy Bengio, Oriol Vinyals, Navdeep Jaitly, and Noam Shazeer, "Scheduled sampling for sequence prediction with recurrent neural networks," in *Proc. NeurIPS*, Montréal, 2015.

[34] Andrey Malinin and Mark Gales, "Predictive uncertainty estimation via prior networks," in *Proc. NeurIPS*, Montréal, 2018.

[35] Andrey Malinin and Mark Gales, "Reverse KL-Divergence training of prior networks: Improved uncertainty and adversarial robustness," in *Proc. NeurIPS*, Vancouver, 2019.

[36] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al., "On the opportunities and risks of foundation models," *arXiv preprint arXiv:2108.07258*, 2021.

[37] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al., "Language models are unsupervised multitask learners," *OpenAI blog*, vol. 1, no. 8, pp. 9, 2019.

[38] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al., "Language models are few-shot learners," in *Proc. NeurIPS*, Virtual, 2020.

[39] Alexander Kolesnikov, Alexey Dosovitskiy, Dirk Weissenborn, Georg Heigold, Jakob Uszkoreit, Lucas Beyer, Matthias Minderer, Mostafa Dehghani, Neil Houlsby, Sylvain Gelly, Thomas Unterthiner, and Xiaohua Zhai, "An image is worth 16x16 words: Transformers for image recognition at scale," in *Proc. ICLR*, Vienna, 2021.

[40] Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever, "Generative pretraining from pixels," in *Proc. ICML*, Virtual, 2020.

[41] Abdelrahman Mohamed, Hung-yi Lee, Lasse Borgholt, Jakob D Havtorn, Joakim Edin, Christian Igel, Katrin Kirchhoff, Shang-Wen Li, Karen Livescu, Lars Maaløe, et al., "Self-supervised speech representation learning: A review," *arXiv preprint arXiv:2205.10643*, 2022.

[42] Jan Chorowski, Grzegorz Ciesielski, Jarosław Dzikowski, Adrian Łańcucki, Ricard Marxer, Mateusz Opala, Piotr Pusz, Paweł Rychlikowski, and Michał Stypułkowski, "Aligned contrastive predictive coding," in *Proc. Interspeech*, Brno, 2021.

[43] Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed, "HuBERT: Self-supervised speech representation learning by masked prediction of hidden units," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 3451–3460, 2021.

[44] Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki Kanda, Takuya Yoshioka, Xiong Xiao, et al., "Wavlm: Large-scale self-supervised pre-training for full stack speech processing," *IEEE Journal of Selected Topics in Signal Processing*, 2022.

[45] Mariana Rodrigues Makiuchi, Kuniaki Uto, and Koichi Shinoda, "Multimodal emotion recognition with high-level speech and text features," in *Proc. ASRU*, Cartagena, 2021.

[46] Edmilson Morais, Ron Hoory, Weizhong Zhu, Itai Gat, Matheus Damasceno, and Hagai Aronowitz, "Speech emotion recognition using self-supervised features," in *Proc. ICASSP*, Toronto, 2022.

[47] Mayank Sharma, "Multi-lingual multi-task speech emotion recognition using wav2vec 2.0," in *Proc. ICASSP*, Toronto, 2022.

[48] Heqing Zou, Yuke Si, Chen Chen, Deepu Rajan, and Eng Siong Chng, "Speech emotion recognition with co-attention based multi-level acoustic information," in *Proc. ICASSP*, Toronto, 2022.

[49] Aaron van den Oord, Yazhe Li, and Oriol Vinyals, "Representation learning with contrastive predictive coding," *arXiv preprint arXiv:1807.03748*, 2018.

[50] Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur, "Librispeech: an asr corpus based on public domain audio books," in *Proc. ICASSP*, South Brisbane, 2015.

[51] Joshua B Tenenbaum and William T. Freeman, "Separating style and content with bilinear models," *Neural computation*, vol. 12, no. 6, pp. 1247–1283, 2000.

[52] Chao Zhang, Zichao Yang, Xiaodong He, and Li Deng, "Multimodal intelligence: Representation learning, information fusion, and applications," *IEEE Journal of Selected Topics in Signal Processing*, vol. 14, no. 3, pp. 478–493, 2020.

[53] Guangzhi Sun, Chao Zhang, and Philip C Woodland, "Combination of deep speaker embeddings for diarisation," *Neural Networks*, vol. 141, pp. 372–384, 2021.

[54] Qiujia Li, Florian L. Kreyssig, Chao Zhang, and Philip C. Woodland, "Discriminative neural clustering for speaker diarisation," in *Proc. SLT*, Shenzhen, 2021.