

Multimodal End-to-End Sparse Model for Emotion Recognition

Wenliang Dai*, Samuel Cahyawijaya*, Zihan Liu, Pascale Fung

Center for Artificial Intelligence Research (CAiRE)

Department of Electronic and Computer Engineering

The Hong Kong University of Science and Technology, Clear Water Bay, Hong Kong

{wdai, scahyawijaya, zliucr}@connect.ust.hk, pascale@ece.ust.hk

Abstract

Existing works on multimodal affective computing tasks, such as emotion recognition, generally adopt a two-phase pipeline, first extracting feature representations for each single modality with hand-crafted algorithms and then performing end-to-end learning with the extracted features. However, the extracted features are fixed and cannot be further fine-tuned on different target tasks, and manually finding feature extraction algorithms does not generalize or scale well to different tasks, which can lead to sub-optimal performance. In this paper, we develop a fully end-to-end model that connects the two phases and optimizes them jointly. In addition, we restructure the current datasets to enable the fully end-to-end training. Furthermore, to reduce the computational overhead brought by the end-to-end model, we introduce a sparse cross-modal attention mechanism for the feature extraction. Experimental results show that our fully end-to-end model significantly surpasses the current state-of-the-art models based on the two-phase pipeline. Moreover, by adding the sparse cross-modal attention, our model can maintain performance with around half the computation in the feature extraction part.

1 Introduction

Humans show their characteristics through not only the words they use, but also the way they speak and their facial expressions. Therefore, in multimodal affective computing tasks, such as emotion recognition, there are usually three modalities: textual, acoustic, and visual. One of the main challenges in these tasks is how to model the interactions between different modalities, as they contain both supplementary and complementary information (Baltrušaitis et al., 2018).

* Equal contribution.

Code is available at: <https://github.com/wenliangdai/Multimodal-End2end-Sparse>

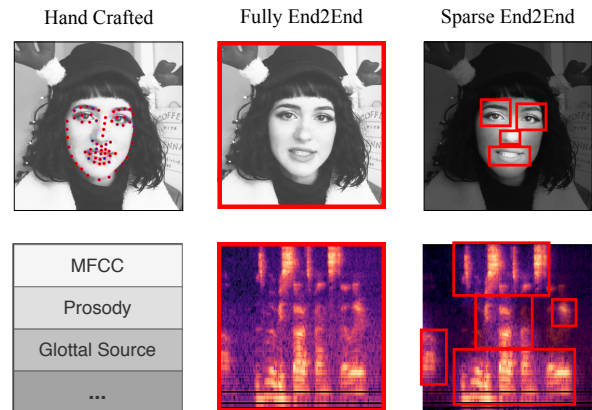


Figure 1: An illustration of feature extraction from hand-crafted model (*left*), fully end-to-end model (*middle*), and sparse end-to-end model (*right*). The red dots represent the keypoints extracted by hand-crafted models. The areas formed by red lines represent the regions of interest that are processed by (sparse) end-to-end models to extract the features.

In the existing works, we discover that a two-phase pipeline is generally used (Zadeh et al., 2018a,b; Tsai et al., 2018, 2019; Rahman et al., 2020). In the first phase, given raw input data, feature representations are extracted with hand-crafted algorithms for each modality separately, while in the second phase, end-to-end multimodal learning is performed using extracted features. However, there are three major defects of this two-phase pipeline: 1) the features are fixed after extraction and cannot be further fine-tuned on target tasks; 2) manually searching for appropriate feature extraction algorithms is needed for different target tasks; and 3) the hand-crafted model considers very few data points to represent higher-level feature, which might not capture all the useful information. These defects can result in sub-optimal performance.

In this paper, we propose a fully end-to-end model that connects the two phases together and optimizes them jointly. In other words, the model receives raw input data and produces the output pre-

dictions, which allows the features to be learned automatically through the end-to-end training. However, the current datasets for multimodal emotion recognition cannot be directly used for the fully end-to-end training, and we thus conduct a data restructuring to make this training possible. The benefits from the end-to-end training are that the features are optimized on specific target tasks, and there is no need to manually select feature extraction algorithms. Despite the advantages of the end-to-end training, it does bring more computational overhead compared to the two-phase pipeline, and exhaustively processing all the data points makes it computationally expensive and prone to overfitting. Thus, to mitigate these side-effects, we also propose a multimodal end-to-end sparse model, a combination of a sparse cross-modal attention mechanism and sparse Convolutional Neural Network (CNN) (Graham and van der Maaten, 2017), to select the most relevant features for the task and reduce the redundant information and noise in the video and audio.

Experimental results show that the simply end-to-end training model is able to consistently outperform the existing state-of-the-art models which are based on the two-phase pipeline. Moreover, the incorporation of the sparse cross-modal attention and sparse CNN is able to greatly reduce the computational cost and maintain the performance.

We summarize our contributions as follows.

- To the best of our knowledge, we are the first to apply a fully end-to-end trainable model for the multimodal emotion recognition task.
- We restructure the existing multimodal emotion recognition datasets to enable the end-to-end training and cross-modal attention based on the raw data.
- We show that the fully end-to-end training significantly outperforms the current state-of-the-art two-phase models, and the proposed sparse model can greatly reduce the computational overhead while maintaining the performance of the end-to-end training. We also conduct a thorough analysis and case study to improve the interpretability of our method.

2 Related Works

Human affect recognition is a popular and widely studied research topic (Mirsamadi et al., 2017;

Zhang and Liu, 2017; Xu et al., 2020; Dai et al., 2020b). In recent years, there is a trend to leverage multimodal information to tackle these research tasks, such as emotion recognition (Busso et al., 2008), sentiment analysis (Zadeh et al., 2016, 2018b), personality trait recognition (Nojavanasghari et al., 2016), etc, have drawn more and more attention. Different methods have been proposed to improve the performance and cross-modal interactions. In earlier works, early fusion (Morency et al., 2011; Pérez-Rosas et al., 2013) and late fusion (Zadeh et al., 2016; Wang et al., 2017) of modalities were widely adopted. Later, more complex approaches were proposed. For example, Zadeh et al. (2017) introduced the Tensor Fusion Network to model the interactions of the three modalities by performing the Cartesian product, while (Wang et al., 2019) used an attention gate to shift the words using the visual and acoustic features. In addition, based on the Transformer (Vaswani et al., 2017), Tsai et al. (2019) introduced the Multimodal Transformer to improve the performance given unaligned multimodal data, and Rahman et al. (2020) introduced a multimodal adaptation gate to integrate visual and acoustic information into a large pre-trained language model. However, unlike some other multimodal tasks (Chen et al., 2017; Yu et al., 2019; Li et al., 2019) using fully end-to-end learning, all of these methods require a feature extraction phase using hand-crafted algorithms (details in Section 5.2), which makes the whole approach a two-phase pipeline.

3 Dataset Reorganization

The fully end-to-end multimodal model requires the inputs to be raw data for the three modalities (visual, textual and acoustic). The existing multimodal emotion recognition datasets cannot be directly applied for the fully end-to-end training for two main reasons. First, the datasets provide split of training, validation and test data for the hand-crafted features as the input of the model and emotion or sentiment labels as the output of the model. However, this dataset split cannot be directly mapped to the raw data since the split indices cannot be matched back to the raw data. Second, the labels of the data samples are aligned with the text modality. However, the visual and acoustic modalities are not aligned with the textual modality in the raw data, which disables the fully end-to-end

training. To make the existing datasets usable for the fully end-to-end training and evaluation, we need to reorganize them according to two steps: 1) align the text, visual and acoustic modalities; 2) split the aligned data into training, validation and test sets.

In this work we reorganize two emotion recognition datasets: Interactive Emotional Dyadic Motion Capture (IEMOCAP) and CMU Multimodal Opinion Sentiment and Emotion Intensity (CMU-MOSEI). Both have multi-class and multi-labelled data for multimodal emotion recognition obtained by generating raw utterance-level data, aligning the three modalities, and creating a new split over the aligned data. In the following section, we will first introduce the existing datasets, and then we will give a detailed description of how we reorganize them.

3.1 IEMOCAP

IEMOCAP (Busso et al., 2008) is a multimodal emotion recognition dataset containing 151 videos. In each video, two professional actors conduct dyadic conversations in English. The dataset is labelled by nine emotion categories, but due to the data imbalance issue, we take the six main categories: *angry*, *happy*, *excited*, *sad*, *frustrated*, and *neutral*. As the dialogues are annotated at the utterance level, we clip the data per utterance from the provided text transcription time, which results in 7,380 data samples in total. Each data sample consists of three modalities: audio data with a sampling rate of 16 kHz, a text transcript, and image frames sampled from the video at 30 Hz. The provided pre-processed data from the existing work (Busso et al., 2008)¹ doesn't provide an identifier for each data sample, which makes it impossible to reproduce it from the raw data. To cope with this problem, we create a new split for the dataset by randomly allocating 70%, 10%, and 20% of data into the training, validation, and testing sets, respectively. The statistics of our dataset split are shown in Table 1.

3.2 CMU-MOSEI

CMU-MOSEI (Zadeh et al., 2018b) comprises 3,837 videos from 1,000 diverse speakers with six emotion categories: *happy*, *sad*, *angry*, *fearful*, *disgusted*, and *surprised*. It is annotated at utterance-level, with a total of 23,259 samples. Each data

¹http://immortal.multicomp.cs.cmu.edu/raw_datasets/processed_data/iemocap

Label	Avg. word length	Avg. clip duration (s)	Train size	Valid size	Test size
Anger	15.96	4.51	757	112	234
Excited	16.79	4.78	736	92	213
Frustrated	17.14	4.71	1298	180	371
Happiness	13.58	4.34	398	62	135
Neutral	13.08	3.90	1214	173	321
Sadness	14.82	5.50	759	118	207

Table 1: Statistics of our IEMOCAP dataset split.

Label	Avg. word length	Avg. clip duration (s)	Train size	Valid size	Test size
Anger	7.75	23.24	3267	318	1015
Disgust	7.57	23.54	2738	273	744
Fear	10.04	28.82	1263	169	371
Happiness	8.14	24.12	7587	945	2220
Sadness	8.12	24.07	4026	509	1066
Surprise	8.40	25.95	1465	197	393

Table 2: Statistics of our CMU-MOSEI dataset split.

sample in CMU-MOSEI consists of three modalities: audio data with a sampling rate of 44.1 kHz, a text transcript, and image frames sampled from the video at 30 Hz. We generate the utterance-level data from the publicly accessible raw CMU-MOSEI dataset.² The generated utterances are perfectly matched with the preprocessed data from the existing work (Zadeh et al., 2018b), but there are two issues with the existing dataset: 1) it includes many misaligned data samples; and 2) many of the samples do not exist in the generated data, and vice versa, in the provided standard split from the CMU MultiModal SDK.³ To cope with the first issue, we perform data cleaning to remove the misaligned samples, which results in 20,477 clips in total. We then create a new dataset split following the CMU-MOSEI split for the sentiment classification task.⁴ The statistics of the new dataset split setting are shown in Table 2.

4 Methodology

4.1 Problem Definition

We define I multimodal data samples as $X = \{(t_i, a_i, v_i)\}_{i=1}^I$, in which t_i is a sequence of words, a_i is a sequence of spectrogram chunks from the

²http://immortal.multicomp.cs.cmu.edu/raw_datasets/processed_data/cmu-mosei/seq_length_20/

³<https://github.com/A2Zadeh/CMU-MultimodalSDK>

⁴http://immortal.multicomp.cs.cmu.edu/raw_datasets/processed_data/cmu-mosei/seq_length_50/mosei_senti_data_noalign.pkl

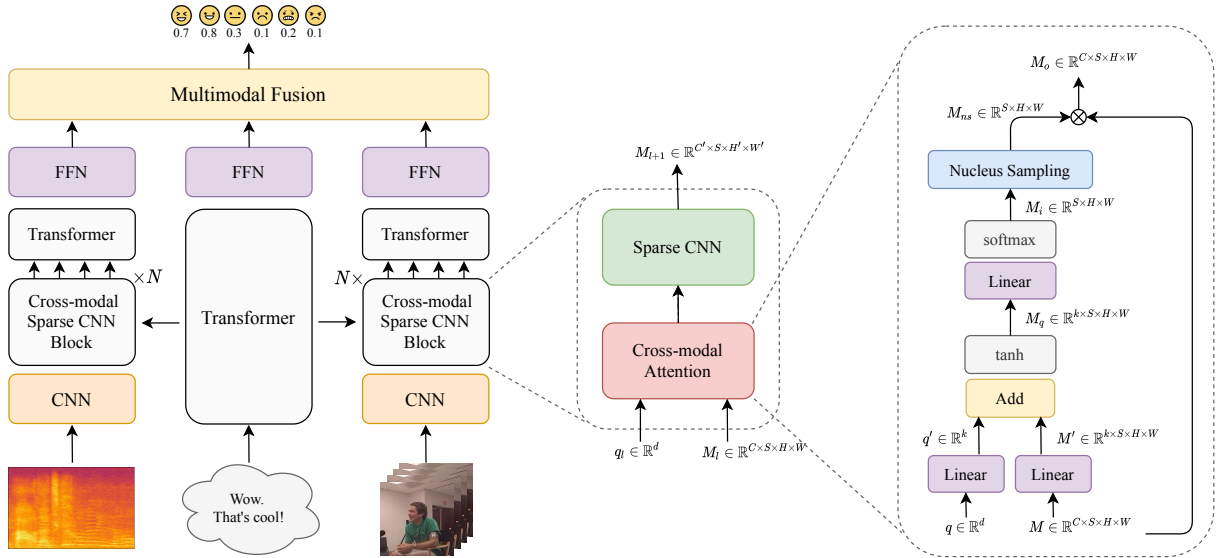


Figure 2: Architecture of our Multimodal End-to-end Sparse Model (MESM). On the left, we show the general architecture flow. In the middle and on the right, we exhibit the details of the cross-modal sparse CNN block, especially the cross-modal attention layer, which is the key to making the CNN model sparse.

audio, and v_i is a sequence of RGB image frames from the video. $Y = \{y_i\}_{i=1}^I$ denotes the annotation for each data sample.

4.2 Fully End-to-End Multimodal Modeling

We build a fully end-to-end model which jointly optimizes the two separate phases (feature extraction and multimodal modelling).

For each spectrogram chunk and image frame in the visual and acoustic modalities, we first use a pre-trained CNN model (an 11-layer VGG (Simonyan and Zisserman, 2014) model) to extract the input features, which are then flattened to vector representations using a linear transformation. After that, we can obtain a sequence of representations for both visual and acoustic modalities. Then, we use a Transformer (Vaswani et al., 2017) model to encode the sequential representations since it contains positional embeddings to model the temporal information. Finally, we take the output vector at the “CLS” token and apply a feed-forward network (FFN) to get the classification scores.

In addition, to reduce GPU memory and align with the two-phase baselines which extract visual features from human faces, we use a MTCNN (Zhang et al., 2016) model to get the location of faces for the image frames before feeding them into the VGG. For the textual modality, the Transformer model is directly used to process the sequence of words. Similar to the visual and acoustic modalities, we consider the feature at the

“CLS” token as the output feature and feed it into a FFN to generate the classification scores. We take a weighted sum of the classification scores from each modality to make the final prediction score.

4.3 Multimodal End-to-end Sparse Model

Although the fully end-to-end model has many advantages over the two-phase pipeline, it also brings much computational overhead. To reduce this overhead without downgrading the performance, we introduce our Multimodal End-to-end Sparse Model (MESM). Figure 2 shows the overall architecture of MESM. In contrast to the fully end-to-end model, we replace the original CNN layers (except the first one for low-level feature capturing) with N cross-modal sparse CNN blocks. A cross-modal sparse CNN block consists of two parts, a cross-modal attention layer and a sparse CNN model that contains two sparse VGG layers and one sparse max-pooling layer.

4.3.1 Cross-modal Attention Layer

The cross-modal attention layer accepts two inputs: a query vector $q \in \mathbb{R}^d$ and a stack of feature maps $M \in \mathbb{R}^{C \times S \times H \times W}$, where C , S , H , and W are the number of channels, sequence length, height, and width, respectively. Then, the cross-modal spatial attention is performed over the feature maps using the query vector. The cross-modal spatial attention

can be formularized in the following steps:

$$M_q = \tanh((W_m M + b_m) \oplus W_q q) \quad (1)$$

$$M_i = \text{softmax}(W_i M_q + b_i) \quad (2)$$

$$M_{ns} = \text{Nucleus Sampling}(M_i) \quad (3)$$

$$M_o = M_{ns} \otimes M, \quad (4)$$

in which $W_m \in \mathbb{R}^{k \times C}$, $W_q \in \mathbb{R}^{k \times d}$, and $W_i \in \mathbb{R}^k$ are linear transformation weights, and $b_m \in \mathbb{R}^k$ and $b_i \in \mathbb{R}^1$ are biases, where k is a pre-defined hyper-parameter, and \oplus represents the broadcast addition operation of a tensor and a vector. In Eq.2, the softmax function is applied to the $(H \times W)$ dimensions, and $M_i \in \mathbb{R}^{S \times H \times W}$ is the tensor of the spatial attention scores corresponding to each feature map. Finally, to make the input feature maps M sparse while reserving important information, firstly, we perform Nucleus Sampling (Holtzman et al., 2019) on M_i to get the top- p portion of the probability mass in each attention score map (p is a pre-defined hyper-parameter in the range of $(0, 1]$). In M_{ns} , the points selected by the Nucleus Sampling are set to one and the others are set to zero. Then, we do broadcast point-wise multiplication between M_{ns} and M to generate the output M_o . Therefore, M_o is a sparse tensor with some positions being zero, and the degree of sparsity is controlled by p .

4.3.2 Sparse CNN

We use the submanifold sparse CNN (Graham and van der Maaten, 2017) after the cross-modal attention layer. It is leveraged for processing low-dimensional data which lies in a space of higher dimensionality. In the multimodal emotion recognition task, we assume that only part of the data is related to the recognition of emotions (an intuitive example is given in Figure 1), which makes it align with the sparse setting. In our model, the sparse CNN layer accepts the output from the cross-modal attention layer, and does convolution computation only at the active positions. Theoretically, in terms of the amount of computation (FLOPs) at a single location, a standard convolution costs $z^2 mn$ FLOPs, and a sparse convolution costs amn FLOPs, where z is the kernel size, m is the number of input channels, n is the number of output channels, and a is the number of active points at this location. Therefore, considering all locations and all layers, the sparse CNN can help to significantly reduce computation.

5 Experiments

5.1 Evaluation Metrics

Following prior works (Tsai et al., 2018; Wang et al., 2019; Tsai et al., 2019; Dai et al., 2020a), we use the accuracy and F1-score to evaluate the models on the IEMOCAP dataset. On the CMU-MOSEI dataset, we use the weighted accuracy instead of the standard accuracy. Additionally, according to Dai et al. (2020a), we use the standard binary F1 rather than the weighted version.

Weighted Accuracy Similar to existing works (Zadeh et al., 2018b; Akhtar et al., 2019), we use the weighted accuracy (WAcc) (Tong et al., 2017) to evaluate the CMU-MOSEI dataset, which contains many more negative samples than positive ones on each emotion category. If normal accuracy is used, a model will still get a fine score when predicting all samples to be negative. The formula of the weighted accuracy is

$$\text{WAcc.} = \frac{TP \times N/P + TN}{2N},$$

in which P means total positive, TP true positive, N total negative, and TN true negative.

5.2 Baselines

For our baselines, we use a two-phase pipeline, which consists of a feature extraction step and an end-to-end learning step.

Feature Extraction We follow the feature extraction procedure in the previous works (Zadeh et al., 2018b; Tsai et al., 2018, 2019; Rahman et al., 2020). For the visual data, we extract 35 facial action units (FAUs) using the OpenFace library⁵ (Baltrušaitis et al., 2015; Baltrušaitis et al., 2018) for the image frames in the video, which capture the movement of facial muscles (Ekman et al., 1980). For the acoustic data, we extract a total of 142 dimension features consisting of 12 dimension bark band energy (BBE) features, 22 dimension mel-frequency cepstral coefficient (MFCC) features, and 108 statistical features from 18 phonological classes. We extract the features per 400 ms time frame using the DisVoice library⁶ (Vásquez-Correa et al., 2018, 2019). For textual data, we use the pre-trained

⁵<https://github.com/TadasBaltrušaitis/OpenFace>

⁶<https://github.com/jcvasquezc/DisVoice>

Model	#FLOPs ($\times 10^9$)	Angry		Excited		Frustrated		Happy		Neutral		Sad		Average	
		Acc.	F1	Acc.	F1	Acc.	F1	Acc.	F1	Acc.	F1	Acc.	F1	Acc.	F1
LF-LSTM	-	71.2	49.4	79.3	57.2	68.2	51.5	67.2	37.6	66.5	47.0	78.2	54.0	71.8	49.5
LF-TRANS	-	81.9	50.7	85.3	57.3	60.5	49.3	85.2	37.6	72.4	49.7	87.4	57.4	78.8	50.3
EmoEmbs [†]	-	65.9	48.9	73.5	58.3	68.5	52.0	69.6	38.3	73.6	48.7	80.8	53.0	72.0	49.8
MuT [†]	-	77.9	60.7	76.9	58.0	72.4	57.0	80.0	46.8	74.9	53.7	83.5	65.4	77.6	56.9
FE2E	8.65	88.7	63.9	89.1	61.9	71.2	57.8	90.0	44.8	79.1	58.4	89.1	65.7	84.5	58.8
MESM ($p = 0.7$)	5.18	88.2	62.8	88.3	61.2	74.9	58.4	89.5	47.3	77.0	52.0	88.6	62.2	84.4	57.4

Table 3: The results on the IEMOCAP dataset. #FLOPs is the number of floating point operations per second. We report the accuracy (Acc.) and the F1-score on six emotion categories: *angry*, *excited*, *frustrated*, *happy*, *neutral* and *sad*. We re-run the models marked by [†], as we use two more categories and the split is different.

Model	#FLOPs ($\times 10^9$)	Angry		Disgusted		Fear		Happy		Sad		Surprised		Average	
		WAcc.	F1	WAcc.	F1	WAcc.	F1	WAcc.	F1	WAcc.	F1	WAcc.	F1	WAcc.	F1
LF-LSTM	-	64.5	47.1	70.5	49.8	61.7	22.2	61.3	73.2	63.4	47.2	57.1	20.6	63.1	43.3
LF-TRANS	-	65.3	47.7	74.4	51.9	62.1	24.0	60.6	72.9	60.1	45.5	62.1	24.2	64.1	44.4
EmoEmbs [†]	-	66.8	49.4	69.6	48.7	63.8	23.4	61.2	71.9	60.5	47.5	63.3	24.0	64.2	44.2
MuT [†]	-	64.9	47.5	71.6	49.3	62.9	25.3	67.2	75.4	64.0	48.3	61.4	25.6	65.4	45.2
FE2E	8.65	67.0	49.6	77.7	57.1	63.8	26.8	65.4	72.6	65.2	49.0	66.7	29.1	67.6	47.4
MESM (0.5)	4.34	66.8	49.3	75.6	56.4	65.8	28.9	64.1	72.3	63.0	46.6	65.7	27.2	66.8	46.8

Table 4: The results on the CMU-MOSEI dataset. WAcc stands for weighted accuracy. We report the accuracy and the F1-score on six emotion categories: *angry*, *disgusted*, *fear*, *happy*, *sad* and *surprised*. We re-run the models marked by [†], as the data we use is unaligned along the sequence length dimension and the split is different.

GloVe (Pennington et al., 2014) word embeddings (glove.840B.300d⁷).

Multimodal Learning As different modalities are unaligned in the data, we cannot compare our method with existing works that can only handle aligned input data. We use four multimodal learning models as baselines: the late fusion LSTM (LF-LSTM) model, the late fusion Transformer (LF-TRANS) model, the Emotion Embeddings (EmoEmbs) model (Dai et al., 2020a), and the Multimodal Transformer (MuT) model (Tsai et al., 2019). They receive the hand-crafted features extracted from the first step as input and give the classification decisions.

5.3 Training Details

We use the Adam optimizer (Kingma and Ba, 2014) for the training of every model we use. For the loss function, we use the binary cross-entropy loss as both of the datasets are multi-class and multi-labelled. In addition, the loss for the positive samples is weighted by the ratio of the number of positive and negative samples to mitigate the imbalance problem. For all of the models, we perform an exhaustive hyper-parameter search to ensure we have solid comparisons. The best hyper-parameters

⁷<https://nlp.stanford.edu/projects/glove/>

are reported in Appendix A. Our experiments are run on an Nvidia 1080Ti GPU, and our code is implemented in the PyTorch (Paszke et al., 2019) framework v1.6.0. We perform preprocessing for the text and audio modalities. For the text modality, we perform word tokenization for our baseline and subword tokenization for our end-to-end model. We limit the length of the text to up to 50 tokens. For the audio modality, we use mel-spectrograms with a window size of 25 ms and stride of 12.5 ms and then chunk the spectrograms per 400 ms time window.

6 Analysis

6.1 Results Analysis

In Table 3, we show the results on the IEMOCAP dataset. Compared to the baselines, the fully end-to-end (FE2E) model surpasses them by a large margin on all the evaluation metrics. Empirically, this shows the superiority of the FE2E model over the two-phase pipeline. Furthermore, our MESM achieves comparable results with the FE2E model, while requiring much less computation in the feature extraction. Here, we only show the results of MESM with the best p value of the Nucleus Sampling. In Section 6.3, we conduct a more detailed discussion of the effects of the top-p values. We further evaluate the methods on the CMU-MOSEI

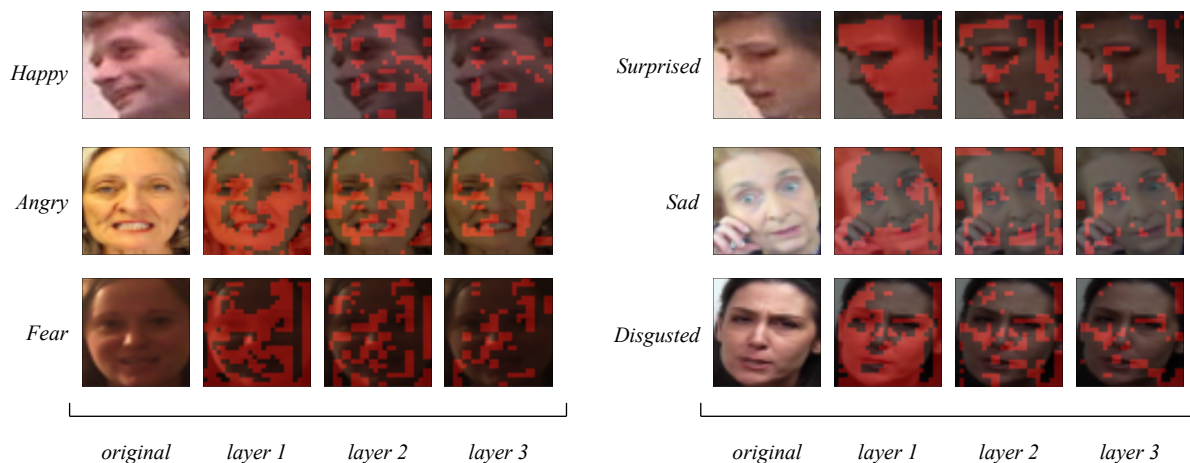


Figure 3: Case study of MESM on six basic emotion categories (happy, sad, angry, surprised, fear, disgusted). From left to right, we show the original image and the Nucleus Sampling ($p = 0.6$) result over points in each attention layer. Red regions represent the points that are computed for the next layer.

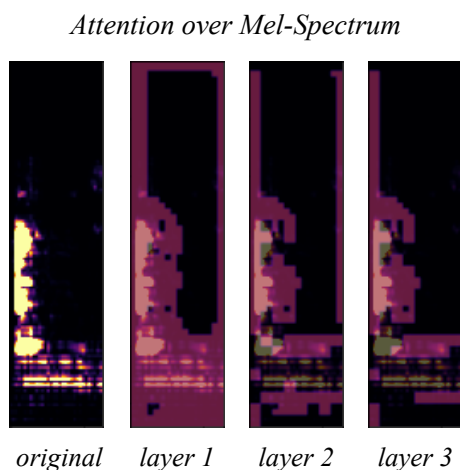


Figure 4: Visualization of cross-modal attention of the acoustic modality. We only show the highest 10% value of mel-spectrogram in each image. From left to right, we show the original image and the Nucleus Sampling ($p = 0.6$) results over points in each attention layer. Red regions represent the active points that will go to the next sparse CNN layer.

dataset and the results are shown in Table 4. We observe similar trends on this dataset.

6.2 Case Study

To improve the interpretability and gain more insights from our model, we visualize the attention maps of our sparse cross-modal attention mechanism on the six basic emotions: happy, sad, angry, surprised, fear, and disgusted. As shown in Figure 3, in general, the models attend to several regions of interest such as the mouth, eyes, eyebrows, and facial muscles between the mouth and

the eyes. We verify our method by comparing the regions that our model captures based on the facial action coding system (FACS) (Ekman, 1997). Following the mapping of FACS to human emotion categories (Basori, 2016; Ahn and Chung, 2017), we conduct empirical analysis to validate the sparse cross-modal attention on each emotion category. For example, the emotion *happy* is highly influenced by raising of the lip on both ends, while *sad* is related to a lowered lip on both ends and downward movement of the eyelids. *Angry* is determined from a narrowed gap between the eyes and thinned lips, while *surprised* is expressed with an open mouth and raising of the eyebrows and eyelids. *Fear* is indicated by a rise of the eyebrows and upper eyelids, and also an open mouth with the ends of the lips slightly moving toward the cheeks. For the emotion *disgusted*, wrinkles near the nose area and movement of the upper lip region are the determinants.

Based on the visualization of the attention maps on the visual data in Figure 3, the MESM can capture most of the specified regions of interest for the six emotion categories. For the emotion *angry*, the sparse cross-modal attention can retrieve the features from the lip region quite well, but it sometimes fails to capture the gap between the eyes. For *surprised*, the eyelids and mouth regions can be successfully captured by MESM, but sometimes the model fails to consider the eyebrow regions. For the acoustic modality, it is hard to analyse the attention in terms of emotion labels. We show a general visualization of the attention maps over the audio data in Figure 4. The model attends to

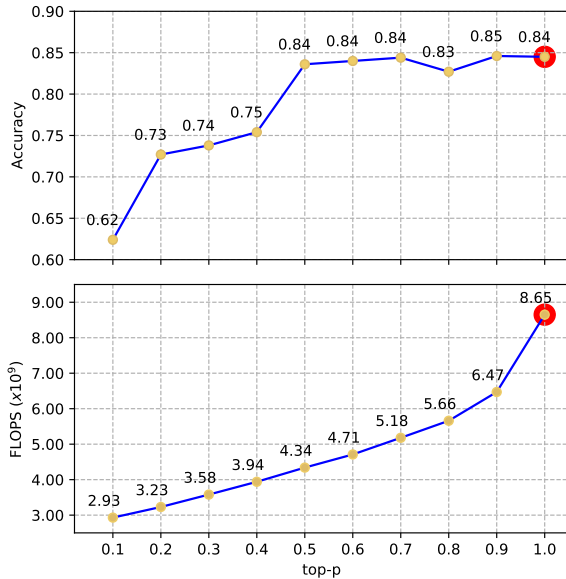


Figure 5: The trend line of the **Top:** Weighted Accuracy and **Bottom:** FLOPs ($\times 10^9$) of the MESM with different top-p values used in the Nucleus Sampling. ● represents performance of MESM, while ● represents performance of the FE2E model

the regions with high spectrum values in the early attention layer, and more points are filtered out after going through further cross-modal attention layers. More visualized examples are provided in Appendix B.

6.3 Effects of Nucleus Sampling

To have an in-depth understanding of the effects of Nucleus Sampling on the MESM, we perform more experiments with different top-p values ranging from 0 to 1, with a step of 0.1. As shown in Figure 5, empirically, the amount of computation is reduced consistently with the decrease of the top-p values. In terms of performance, with a top-p value from 0.9 to 0.5, there is no significant drop in the evaluation performance. Starting from 0.5 to 0.1, we can see a clear downgrade in the performance, which means some of the useful information for recognizing the emotion is excluded. The inflection point of this elbow shaped trend line can be an indicator to help us make a decision on the value of the top-p. Specifically, with a top-p of 0.5, the MESM can achieve comparable performance to the FE2E model with around half of the FLOPs in the feature extraction.

Model	Mods.	Avg. Acc	Avg. F1
FE2E	TAV	84.5	58.5
	TA	83.7	54.0
	TV	82.8	55.7
	VA	81.2	54.4
	T	80.8	50.0
	A	73.3	44.9
	V	78.2	49.8
MESM	TAV	84.4	57.3
	TA	83.6	56.7
	TV	82.1	56.0

Table 5: Results of the ablation study of our fully end-to-end model (FE2E) and multimodal end-to-end sparse model (MESM) on the IEMOCAP dataset. In the Mods. (modalities) column, the T/A/V indicates the existence of the textual (T), acoustic (A), and visual (V) modalities.

7 Ablation Study

We conduct a comprehensive ablation study to further investigate how the models perform when one or more modalities are absent. The results are shown in Table 5. Firstly, we observe that the more modalities the more improvement in the performance. TAV, representing the presence of all three modalities, results in the best performance for both models, which shows the effectiveness of having more modalities. Secondly, with only a single modality, the textual modality results in better performance than the other two, which is similar to the results of previous multimodal works. This phenomenon further validates that using textual (T) to attend to acoustic (A) and visual (V) in our cross-modal attention mechanism is a reasonable choice. Finally, with two modalities, the MESM can still achieve a performance that is on par with the FE2E model or is even slightly better.

8 Conclusion and Future Work

In this paper, we first compare and contrast the two-phase pipeline and the fully end-to-end (FE2E) modelling of the multimodal emotion recognition task. Then, we propose our novel multimodal end-to-end sparse model (MESM) to reduce the computational overhead brought by the fully end-to-end model. Additionally, we reorganize two existing datasets to enable fully end-to-end training. The empirical results demonstrate that the FE2E model has an advantage in feature learning

and surpasses the current state-of-the-art models that are based on the two-phase pipeline. Furthermore, MESM is able to halve the amount of computation in the feature extraction part compared to FE2E, while maintaining its performance. In our case study, we provide a visualization of the cross-modal attention maps on both visual and acoustic data. It shows that our method can be interpretable, and the cross-modal attention can successfully select important feature points based on different emotion categories. For future work, we believe that incorporating more modalities into the sparse cross-modal attention mechanism is worth exploring since it could potentially enhance the robustness of the sparsity (selection of features).

Acknowledgement

This work is funded by MRP/055/18 of the Innovation Technology Commission, the Hong Kong SAR Government.

References

- Duck-ki Ahn and Jean-Hun Chung. 2017. [A study on character's emotional appearance in distinction focused on 3d animation "inside out"](#). *Journal of Digital Convergence*, 15:361–368.
- Md Shad Akhtar, Dushyant Singh Chauhan, Deepanway Ghosal, Soujanya Poria, Asif Ekbal, and Pushpak Bhattacharyya. 2019. Multi-task learning for multi-modal emotion recognition and sentiment analysis. *arXiv preprint arXiv:1905.05812*.
- Tadas Baltrušaitis, Chaitanya Ahuja, and Louis-Philippe Morency. 2018. Multimodal machine learning: A survey and taxonomy. *IEEE transactions on pattern analysis and machine intelligence*, 41(2):423–443.
- Tadas Baltrušaitis, Amir Zadeh, Yao Chong Lim, and Louis-Philippe Morency. 2018. Openface 2.0: Facial behavior analysis toolkit. In *2018 13th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2018)*, pages 59–66. IEEE.
- T. Baltrušaitis, M. Mahmoud, and P. Robinson. 2015. [Cross-dataset learning and person-specific normalization for automatic action unit detection](#). In *2015 11th IEEE International Conference and Workshops on Automatic Face and Gesture Recognition (FG)*, volume 06, pages 1–6.
- Ahmad Hoirul Basori. 2016. [Emotional facial expression based on action units and facial muscle](#). *International Journal of Electrical and Computer Engineering (IJECE)*, 6:2478–2487.
- Carlos Busso, Murtaza Bulut, Chi-Chun Lee, Abe Kazemzadeh, Emily Mower, Samuel Kim, Jeanette N Chang, Sungbok Lee, and Shrikanth S Narayanan. 2008. Iemocap: Interactive emotional dyadic motion capture database. *Language resources and evaluation*, 42(4):335.
- Long Chen, Hanwang Zhang, Jun Xiao, L. Nie, Jian Shao, Wei Liu, and Tat-Seng Chua. 2017. Sca-cnn: Spatial and channel-wise attention in convolutional networks for image captioning. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6298–6306.
- Wenliang Dai, Zihan Liu, Tiezheng Yu, and Pascale Fung. 2020a. Modality-transferable emotion embeddings for low-resource multimodal emotion recognition. *ArXiv*, abs/2009.09629.
- Wenliang Dai, Tiezheng Yu, Zihan Liu, and Pascale Fung. 2020b. Kungfupanda at semeval-2020 task 12: Bert-based multi-task learning for offensive language detection. *ArXiv*, abs/2004.13432.
- Paul Ekman, Wallace V Freisen, and Sonia Ancoli. 1980. Facial signs of emotional experience. *Journal of personality and social psychology*, 39(6):1125.
- Rosenberg Ekman. 1997. *What the face reveals: Basic and applied studies of spontaneous expression using the Facial Action Coding System (FACS)*. Oxford University Press, USA.
- Benjamin Graham and Laurens van der Maaten. 2017. Submanifold sparse convolutional networks. *arXiv preprint arXiv:1706.01307*.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2019. The curious case of neural text degeneration. *arXiv preprint arXiv:1904.09751*.
- Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Liunian Harold Li, Mark Yatskar, Da Yin, C. Hsieh, and Kai-Wei Chang. 2019. Visualbert: A simple and performant baseline for vision and language. *ArXiv*, abs/1908.03557.
- Seyedmahdad Mirsamadi, E. Barsoum, and C. Zhang. 2017. Automatic speech emotion recognition using recurrent neural networks with local attention. *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2227–2231.
- Louis-Philippe Morency, Rada Mihalcea, and Payal Doshi. 2011. Towards multimodal sentiment analysis: Harvesting opinions from the web. In *Proceedings of the 13th international conference on multimodal interfaces*, pages 169–176.
- Behnaz Nojavanasghari, Deepak Gopinath, Jayanth Koushik, Tadas Baltrušaitis, and Louis-Philippe

- Morency. 2016. Deep multimodal fusion for persuasiveness prediction. In *Proceedings of the 18th ACM International Conference on Multimodal Interaction*, pages 284–288.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. In *Advances in neural information processing systems*, pages 8026–8037.
- Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543.
- Verónica Pérez-Rosas, Rada Mihalcea, and Louis-Philippe Morency. 2013. Utterance-level multimodal sentiment analysis. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 973–982.
- Wasifur Rahman, Md Kamrul Hasan, Sangwu Lee, AmirAli Bagher Zadeh, Chengfeng Mao, Louis-Philippe Morency, and Ehsan Hoque. 2020. Integrating multimodal information in large pretrained transformers. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2359–2369.
- Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*.
- Edmund Tong, Amir Zadeh, Cara Jones, and Louis-Philippe Morency. 2017. Combating human trafficking with multimodal deep models. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1547–1556.
- Yao-Hung Hubert Tsai, Shaojie Bai, Paul Pu Liang, J Zico Kolter, Louis-Philippe Morency, and Ruslan Salakhutdinov. 2019. Multimodal transformer for unaligned multimodal language sequences. In *Proceedings of the conference. Association for Computational Linguistics. Meeting*, volume 2019, page 6558. NIH Public Access.
- Yao-Hung Hubert Tsai, Paul Pu Liang, Amir Zadeh, Louis-Philippe Morency, and Ruslan Salakhutdinov. 2018. Learning factorized multimodal representations. *arXiv preprint arXiv:1806.06176*.
- J. C. Vásquez-Correa, P. Klumpp, J. R. Orozco-Arroyave, and E. Nöth. 2019. Phonet: A tool based on gated recurrent neural networks to extract phonological posteriors from speech. In *INTERSPEECH*.
- Juan Camilo Vásquez-Correa, J. R. Orozco-Arroyave, T. Bocklet, and E. Nöth. 2018. Towards an automatic evaluation of the dysarthria level of patients with parkinson’s disease. *Journal of communication disorders*, 76:21–36.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008.
- Haohan Wang, Aaksha Meghawat, Louis-Philippe Morency, and Eric P Xing. 2017. Select-additive learning: Improving generalization in multimodal sentiment analysis. In *2017 IEEE International Conference on Multimedia and Expo (ICME)*, pages 949–954. IEEE.
- Yansen Wang, Ying Shen, Zhun Liu, Paul Pu Liang, Amir Zadeh, and Louis-Philippe Morency. 2019. Words can shift: Dynamically adjusting word representations using nonverbal behaviors. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 7216–7223.
- Peng Xu, Zihan Liu, Genta Indra Winata, Zhaojiang Lin, and Pascale Fung. 2020. Emograph: Capturing emotion correlations using graph networks. *ArXiv*, abs/2008.09378.
- Zhou Yu, J. Yu, Yuhao Cui, D. Tao, and Qi Tian. 2019. Deep modular co-attention networks for visual question answering. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6274–6283.
- Amir Zadeh, Minghai Chen, Soujanya Poria, Erik Cambria, and Louis-Philippe Morency. 2017. Tensor fusion network for multimodal sentiment analysis. *arXiv preprint arXiv:1707.07250*.
- Amir Zadeh, Paul Pu Liang, Navonil Mazumder, Soujanya Poria, Erik Cambria, and Louis-Philippe Morency. 2018a. Memory fusion network for multi-view sequential learning. *arXiv preprint arXiv:1802.00927*.
- Amir Zadeh, Rowan Zellers, Eli Pincus, and Louis-Philippe Morency. 2016. Mosei: multimodal corpus of sentiment intensity and subjectivity analysis in online opinion videos. *arXiv preprint arXiv:1606.06259*.
- AmirAli Bagher Zadeh, Paul Pu Liang, Soujanya Poria, Erik Cambria, and Louis-Philippe Morency. 2018b. Multimodal language analysis in the wild: Cmu-mosei dataset and interpretable dynamic fusion graph. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2236–2246.
- Kaipeng Zhang, Zhanpeng Zhang, Zhifeng Li, and Yu Qiao. 2016. Joint face detection and alignment using multitask cascaded convolutional networks. *IEEE Signal Processing Letters*, 23(10):1499–1503.
- L. Zhang and B. Liu. 2017. Sentiment analysis and opinion mining. In *Encyclopedia of Machine Learning and Data Mining*.

A Hyper-parameter Settings

	IEMOCAP		CMU-MOSEI	
	FE2E	MESM	FE2E	MESM
Batch size	8	8	8	8
Learning rate	5e-5	5e-5	5e-5	5e-5
Dim	64	64	64	64
#Heads	4	4	4	4
#Layers	4	4	4	4
Max text len	50	100	50	100
N	-	3	-	3

Table 6: The best hyper-parameters used in training for the two datasets.

B Case Study on Acoustic Modality

We provide more visualized examples of the sparse cross-modal attention maps of the acoustic modality in Figure 6.

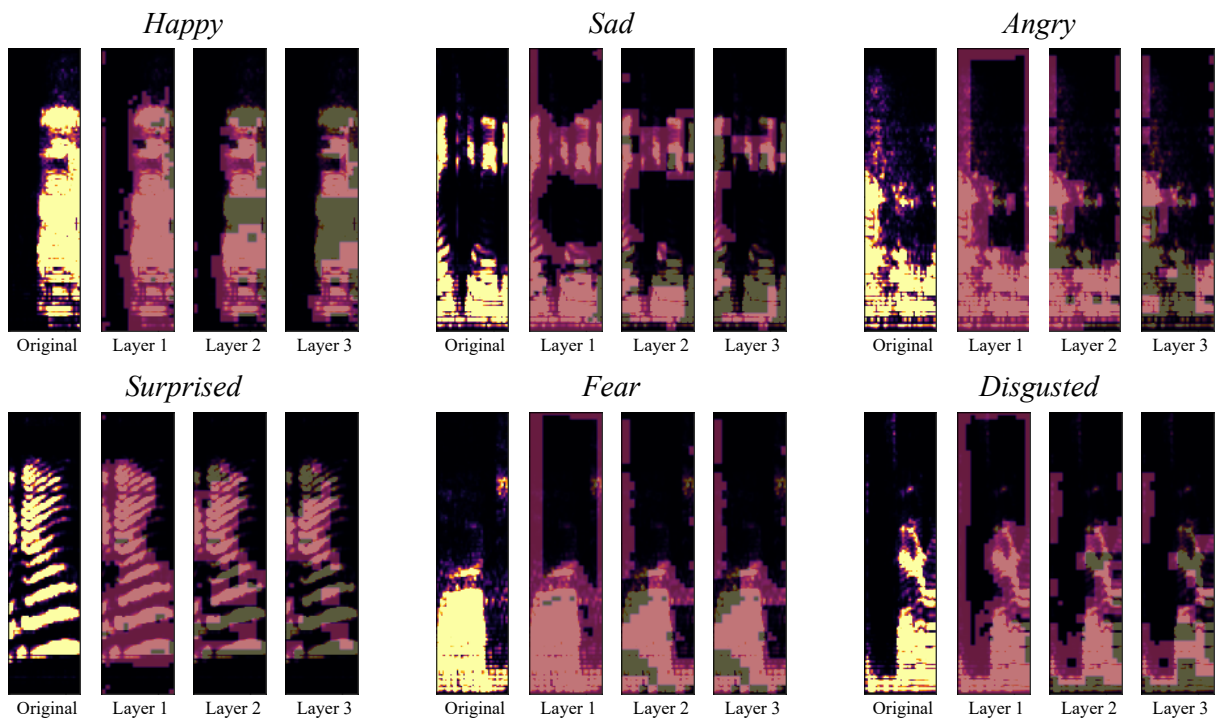


Figure 6: Case study of the sparse cross-modal attention maps on six basic emotion categories (*happy*, *sad*, *angry*, *surprised*, *fear*, *disgusted*) on the audio modality. From the left to right, we show the original image and the Nucleus Sampling results over feature points in each attention layer. Red regions represent the active points that will be computed in the next sparse layer.