# ON THE USE OF MODALITY-SPECIFIC LARGE-SCALE PRE-TRAINED ENCODERS FOR MULTIMODAL SENTIMENT ANALYSIS

*Atsushi Ando, Ryo Masumura, Akihiko Takashima, Satoshi Suzuki, Naoki Makishima,*
*Keita Suzuki, Takafumi Moriya, Takanori Ashihara, Hiroshi Sato*

NTT Corporation

## ABSTRACT

This paper investigates the effectiveness and implementation of modality-specific large-scale pre-trained encoders for multimodal sentiment analysis (MSA). Although the effectiveness of pre-trained encoders in various fields has been reported, conventional MSA methods employ them for only linguistic modality, and their application has not been investigated. This paper compares the features yielded by large-scale pre-trained encoders with conventional heuristic features. One each of the largest pre-trained encoders publicly available for each modality are used; CLIP-ViT, WavLM, and BERT for visual, acoustic, and linguistic modalities, respectively. Experiments on two datasets reveal that methods with domain-specific pre-trained encoders attain better performance than those with conventional features in both unimodal and multimodal scenarios. We also find it better to use the outputs of the intermediate layers of the encoders than those of the output layer. The codes are available at `https://github.com/ando-hub/MSA_Pretrain`.

***Index Terms***— Multimodal Sentiment Analysis, Large-Scale Pre-trained Encoder

## 1. INTRODUCTION

Multimodal sentiment analysis (MSA) is the technology to estimate the sentiment of a target speaker from multimodal information such as visual, acoustic, and linguistic modalities. Since sentimental cues appear in various aspects such as facial expression, tone, and phrases, MSA performs better than alternatives using a single modality such as facial expression recognition [1], speech emotion recognition [2], and sentiment analysis from text [3].

Most conventional studies have focused on modeling the interactions of multiple modalities and modality-specific information. Some aim to model the interactions of short-term features of each modality [4, 5]. They can use local characteristics across modalities, e.g., facial expression changes or prosody during a particular word, for enhancing prediction performance. Others extract sequence-level representations of the speaker's sentiment in the individual modalities,

then estimate the sentiment level from all of the sequence-level representations [6, 7]. These representations have been evaluated against each other to learn modality-invariant and modality-specific information, with the goal being to improve robustness against missing information of specific modalities such as facial occlusion [8, 9]. The studies employ heuristic features and/or the prediction results of the model such as head pose, gaze, and facial landmarks.

Recently, several MSA studies have utilized large-scale pre-trained encoders in addition to developing model structures [10–13]. The pre-trained encoder is a part of the model trained in other tasks known as upstream tasks. The advantage of the pre-trained encoder is that it enables the transfer of common knowledge of upstream tasks, which yields better cues for a target downstream task compared to training from scratch. It has also been reported that a larger pre-trained encoder trained on a large amount of upstream data offers better performance in downstream tasks [14, 15]. Large-scale pre-training encoders have significantly enhanced various downstream tasks in visual, acoustic, and linguistic modalities [16–18]. Though the introduction of the pre-trained encoder has improved MSA performance, the conventional studies suffer from two omissions. First, they employ the large-scale pre-trained encoder only in linguistic modality, not in visual and acoustic modalities. Second, how to apply the pre-trained encoders has not been investigated. Some studies use the weighted sum results of the hidden states extracted from each layer of the pre-trained encoder for speech and speaker recognition [17] since it is empirically known that different knowledge is extracted in the different layers in the pre-trained encoder [19, 20]. However, the previous work employed only the output of the final layer of the pre-trained model, which may be less effective for MSA than the use of the hidden states of the intermediate layers.

This paper investigates the following two research questions: (i) Are the features based on the modality-specific pre-trained encoder more effective than the conventional features in multimodal, and even unimodal, scenarios? (ii) How to apply the modality-specific pre-trained encoders? Three large-scale pre-trained encoders that are currently available for each modality are examined; CLIP Vision Transformer (ViT) [16],

WavLM [17], and BERT [18] for visual, acoustic, and linguistic modalities, respectively. This paper introduces a simple sequence-level cross-modal model, Unimodal Encoders and Gated Decoder (UEGD), which enables comparison of multimodal and unimodal performances on the same model structure. We evaluate three types of representations from the pre-trained encoders that have been used in other downstream tasks; the output, each of the intermediate outputs, and a weighted sum of the intermediate outputs. Experiments on two public datasets, CMU-MOSI [21] and CMU-MOSEI [22], reveal the answers to our research questions: (i) Large-scale pre-trained encoders improve sentiment analysis performance in both unimodal and multimodal scenarios. The UEGD model with pre-trained encoders achieves state-of-the-art performance in regression tasks on CMU-MOSEI. It is also found that the pre-trained encoder is particularly effective in the acoustic modality. (ii) Using one of the late middle intermediate layers of the pre-trained encoder yield better performance than the final layer output and the weighted sum of the intermediate layer outputs.

## 2. RELATED WORK

### 2.1. Multimodal Sentiment Analysis

The conventional MSA methods can be categorized into two approaches; early-fusion and late-fusion.

The early-fusion approach aims to capture interactive information from low-level features such as frame-level features in visual and acoustic modalities and word-level features in linguistic modality. Multimodal Transformer (MulT) employs multiple crossmodal transformers to capture bi-modal interactions [5]. Multimodal Adaptation Gate (MAG) focuses on integrating linguistic features with other modalities, visual and acoustic factors, by using a gating structure in addition to the cross-modal self-attention [10, 11]. The advantage of this approach lies in capturing local characteristics across modalities. However, it can only be used when two or more modalities are available.

Late-fusion integrates utterance-level representations of modalities to predict sentiment. Tensor Fusion Network (TFN) explicitly models uni-, bi-, and tri-modal interactions as outer products of utterance-level embeddings [6]. Modality-Invariant and -Specific Representations (MISA) extract modality-invariant/-specific utterance-level representations by introducing similarity and difference losses of representations [12]. Self-Supervised Multi-task Multimodal sentiment analysis (Self-MM) jointly learns multimodal and unimodal subtasks from utterance-level representations to supplement modality-specific information [13].

The MSA model in this paper uses a late-fusion approach since it enables performance comparisons in unimodal and multimodal scenarios on the same model structure.

### 2.2. Large-Scale Pre-Trained Encoders

Large-scale pre-trained encoders have received significant attention in recent years. In this framework, an upstream model is trained by a large amount of training data, then a small amount of labeled data is used to adapt downstream tasks in combination with a part of the upstream model. The tasks that do not require human annotation, such as contrastive learning [23] or self-supervised learning [15], are used as the upstream task. It has been reported that larger models trained by large amounts of upstream data show higher performance in many downstream tasks. Pre-trained encoders comprising a stack of multiple transformer encoders [24] are often used. This approach was first used with great success in natural language processing [18, 25] and is now widely used in computer vision [16] and speech processing [17, 19]. It has been empirically reported that low-level features are extracted in the layers closer to the input, e.g., phrase-level information in linguistic encoder, and high-level features are obtained in the layers closer to the output like long-distance dependency information [17, 20].

This work employs the pre-trained models that are some of the largest and publicly available ones; CLIP-ViT [16], WavLM [17], and BERT [18] in visual, acoustic, and linguistic modalities, respectively. CLIP-ViT extracts an image embedding from a single image, while WavLM and BERT yield sequence-level embeddings from a series of an audio waveform and word tokens.

## 3. MULTIMODAL SENTIMENT ANALYSIS BY MODALITY-SPECIFIC PRE-TRAINED ENCODERS

### 3.1. Model Structure

Let $\boldsymbol{S}_v, \boldsymbol{S}_a, \boldsymbol{S}_l$ be input data of visual, acoustic, and linguistic modalities, respectively, and $y$ be the corresponding sentiment value of the utterance. Multimodal sentiment analysis is defined as the regression task of determining $y$ from $\boldsymbol{S}_v, \boldsymbol{S}_a, \boldsymbol{S}_l$,
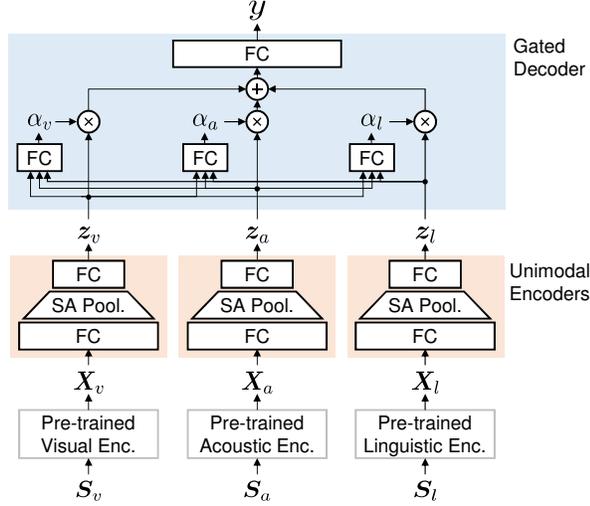
$$\hat{y} = f\left(\boldsymbol{S}_v, \boldsymbol{S}_a, \boldsymbol{S}_l; \Theta\right), \tag{1}$$

where $\hat{y}$ is the predicted sentiment value, $f(\cdot)$ is the regression function determined by the regression model, and $\Theta$ is a parameter set of the regression model.

This paper uses Unimodal Encoders and Gated Decoder (UEGD) as the regression model, see Fig. 1. The inputs of the UEGD model are modal-dependent low-level features $\boldsymbol{X}_v, \boldsymbol{X}_a, \boldsymbol{X}_l$ extracted from $\boldsymbol{S}_v, \boldsymbol{S}_a, \boldsymbol{S}_l$ with pre-trained encoders, see Section 3.2 for details. UEGD extracts utterance-level embeddings of each modality $\boldsymbol{z}_v, \boldsymbol{z}_a, \boldsymbol{z}_l$,

$$\boldsymbol{z}_m = \mathsf{UnimodalEncoder}\left(\boldsymbol{X}_m; \theta_m^{(e)}\right), \tag{2}$$

where $m \in \{v, a, l\}$ and $\mathsf{UnimodalEncoder}(\cdot)$ is a projection function from low-level features to the utterance-level

**Fig. 1**: Overview of the Unimodal Encoders and Gated Decoder (UEGD) model.



**Fig. 2**: The weighted sum of the pre-trained encoder outputs.

**Table 1**: The numbers of clips in the datasets.

| | # Train | # Valid. | # Test | # Total |
|---|---|---|---|---|
| CMU-MOSI | 1284 | 229 | 686 | 2199 |
| CMU-MOSEI | 16326 | 1871 | 4659 | 22856 |

embedding in each modality. $\theta_m^{(e)}$ is a set of parameters of the unimodal encoder. $z_v, z_a, z_l \in \mathbb{R}^F$ and $F$ are dimensions of the embeddings. The unimodal encoder consists of Fully-Connected (FC) and Self-Attentive (SA) pooling layers. SA pooling layer allows hidden vectors having specific intervals in the utterance-level features to be strongly reflected, which is desirable in sentiment analyses since sentiment cues will appear in limited regions of input data.

The utterance-level embeddings are integrated and projected by the gated decoder to yield the predicted result,

$$\hat{y} = \mathsf{GatedDecoder}\left(z_v, z_a, z_l; \theta^{(d)}\right), \quad (3)$$

where $\mathsf{GatedDecoder}(\cdot)$ is a function projecting the estimated sentiment values from the embeddings. $\theta^{(d)}$ is a set of parameters of the gated decoder. The gated decoder has a gate layer that integrates the embeddings by the weighted sum,
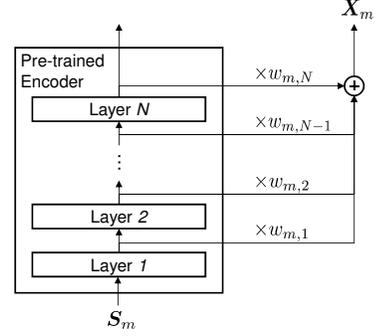
$$\alpha_m = \sigma\left(W_m\left[z_v; z_a; z_l\right] + b_m\right), \quad (4)$$

$$\bar{z} = \sum_{m \in \{v, a, l\}} \alpha_m z_m, \quad (5)$$

where $\alpha_m$ represents the gate weight of the specific modality $m$ calculated by the supervector concatenated by embeddings of all modalities. $\{W_m, b_m\} \in \theta^{(d)}$ is the set of weight and bias parameters. $\sigma(\cdot)$ is a sigmoid function that constraints the gate weight to lie between 0 to 1. The integrated embedding $\bar{z}$ is input to the FC layers to obtain $\hat{y}$. One of the advantages of the gated decoder is that it allows quantifying the contribution of each modality to the prediction result as the gate weight[1].

The model parameters $\Theta = \{\theta_v^{(e)}, \theta_a^{(e)}, \theta_l^{(e)}, \theta^{(d)}\}$ are optimized by L1 loss of ground truth $y$ and the predicted result

[1]In the preliminary experiments, the gate fusion performed the same as the concatenation of the utterance-level embeddings.

$\hat{y}$. Note that all the modality-specific pre-trained encoders are frozen during training since they have so many parameters that making them trainable would lead to overfitting.

Another advantage of the UEGD model is that it can be used for unimodal as well as multimodal cases on the same structure. In the case of unimodal input, the utterance-level embeddings of the unused modalities $z_m$ are set to $0 \in \mathbb{R}^F$. Experiments have confirmed that the gate weights of the unused modalities approach zero, which means the UEGD model evaluates sentiment from specific unimodal inputs.

### 3.2. Feature Extraction with the Pre-Trained Encoder

Three types of the extracted features $X_m$ yielded by each of the modality-specific pre-trained encoders are investigated; the output, any intermediate output, and a weighted sum of the intermediate outputs. The first two represent the output vectors of the final layer or the intermediate layers of the pre-trained encoder, respectively. The weighted sum uses the weighted sum of the hidden states extracted from each layer of the pre-trained model as reported in the conventional work of [17], see Fig. 2. The weights of the intermediate layers are parameters that are optimized simultaneously with the parameters of the UEGD model by the labeled data.

## 4. EXPERIMENTS

### 4.1. Datasets

Two public multimodal datasets were used in the experiments; CMU-MOSI [21] and CMU-MOSEI [22]. Both contain short clips from YouTube videos with sentiment labels given by human annotators. Each clip contains only one person and an utterance of a few seconds. The label takes values from -3 to

+3. The total numbers of videos are 98 and 3178 on CMU-MOSI and CMU-MOSEI, respectively, which almost match the numbers of speakers. The datasets were divided into training, validation, and test subsets as in the conventional studies [12, 13]. The numbers of clips are shown in Table 1.

## 4.2. Setup

The pre-trained encoders of the individual modalities were as follows. *CLIP ViT-L/14*[2], *WavLM Large*[3], and *BERT-Large-uncased*[4] were used in visual, acoustic, and linguistic modalities, respectively. The training data consisted of 400 M image-text pairs collected from the Internet for CLIP ViT, 94 k hours of speech from audiobooks, podcasts and meetings for WavLM, 3.3 B words from Wikipedia and books. The numbers of hidden layers and embedding sizes were, for all encoders, 24 and 1024, respectively. Face detection by MTCNN [26] was applied to get a $256\times256$-pixel face image before extracting visual embeddings. For visual modality, the outputs of the `[class]` token position were used as the pre-trained encoder outputs for each frame. Face detection and visual feature extraction were applied at 3 fps as in conventional visual feature extraction [6, 22]. The frame shift length of the audio features yielded by the encoders was originally 20 ms and subsampled to 1/10 to get 5 fps features.

Two sets of the conventional features were used: those the same as [27][5] and [13][6]. The former is composed of 35-dimensional FACET facial expression analysis results as visual, 74-dimensional COVEREP hand-crafted features as acoustic, and 300-dimensional GloVe embeddings as linguistic features. The latter have the same visual and acoustic features and 768-dimensional BERT-Base output embeddings as linguistic features. As in the previous work, only 20- and 5-dimensional features were used as the visual and acoustic features in CMU-MOSI.

The compared features from the domain-specific pre-trained encoders were the output (Enc. output), the combinations of the single best of the intermediate outputs (Enc. mid-best), and the weighted sum of the intermediate outputs (Enc. weighted). According to the correlation coefficient of the validation set in the unimodal scenario described in Section 4.3.2, we used the 15 th, 21 st, and 20 th intermediate layers of visual, acoustic, and linguistic encoders respectively as Enc. mid-best in CMU-MOSI, while 19 th, 21 st, and 21 st in CMU-MOSEI.

The hyper-parameters of the proposed UEGD model were as follows. The uni-modal encoder is composed of a fully-connected layer with 256 hidden units, self-attentive pooling with 4 heads, and a fully-connected layer with 128 hidden units, i.e., the size of the utterance-level embedding was 128.

---

[2] https://github.com/openai/CLIP
[3] https://github.com/microsoft/unilm/tree/master/wavlm
[4] https://pypi.org/project/pytorch-pretrained-bert/
[5] https://github.com/A2Zadeh/CMU-MultimodalSDK
[6] https://github.com/thuiar/MMSA/tree/master/src/MMSA

The gated decoder consisted of a gating layer and two fully-connected layers with 128 and 1 hidden units, respectively. We used the above hidden units in CMU-MOSEI and half of them in CMU-MOSI because the amount of the labeled training data was limited in CMU-MOSI. Layer normalization and ReLU activation functions were applied after all fully-connected layers except for the output of the uni-modal encoder and the multimodal decoder. The dropout rate was 0.2. Batchsize was 16. We used the Adam optimizer, and the learning rate was 0.0001 with warmup and cosine annealing. Early stopping was applied via average loss of the validation set. In the training step, masking was applied to a maximum of 20% of the time and feature dimensions as in SpecAugment [28] to prevent overfitting. We employed the same evaluation metrics as the conventional studies: Mean Absolute Error (MAE), pearson correlation coefficient (Corr), Accuracy (Acc), and Weighted F1 score (F-score). The last two were the results of two-class classification tasks that predict negative/non-negative or negative/positive (exclude zero). In all experiments, we ran trials five times; the average performance is taken as the final result.

## 4.3. Results

### 4.3.1. Evaluations of the Multimodal Model

Prediction performances by the conventional methods and the UEGD models with the conventional and pre-trained encoder output-based features are shown in Table 2. Compared to the UEGD models, the encoder-based features showed better performances in MAE and Corr on both datasets. Furthermore, all the encoder-based results with the UEGD models achieved better MAEs and correlation results than the conventional methods on CMU-MOSEI, even though they were based on a simple model and loss function. On the other hand, the proposed method was inferior to several conventional methods on CMU-MOSI. This indicates that the proposed UEGD model may not be optimized for small training data. For example, MAG [10] is explicitly designed to use the linguistic features mainly for prediction, while the proposed method has to learn such knowledge from the training data, which is difficult if the training data is limited. From these results, we consider that pre-trained encoders are more effective for multimodal sentiment analysis than the conventional features, especially in the case of a large amount of labeled training data.

The contributions of the individual modalities were also compared for both the conventional and the encoder-based features from the gate weights of the UEGD model. The distributions of the gate weights in CMU-MOSEI are shown in Fig. 3. The gate weights of the acoustic modality using pre-trained outputs were distributed at higher values than those yielded by the conventional features. This indicates that there was an improvement in the acoustic modality features, which results in better MSA performance.

**Table 2**: Performance of the multimodal models. In the Feature column, Conv and Conv-BERT are sets of the conventional features including GloVe and BERT-Base features from linguistic modality, respectively. In Acc-2 and F-Score, the left of "/" is "negative/non-negative" score and the right is "negative/positive" performance. Results of $^{\dagger}$, $^{\ddagger}$ and $\diamond$ are from [12], [10], and [13], respectively.

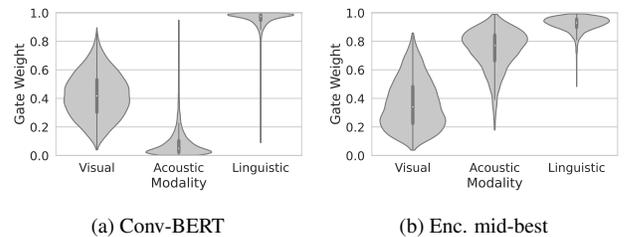| | Feature | CMU-MOSI | | | | CMU-MOSEI | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | MAE↓ | Corr↑ | Acc-2↑ | F-Score↑ | MAE↓ | Corr↑ | Acc-2↑ | F-Score↑ |
| TFN [6]$^{\dagger}$ | Conv | 0.970 | 0.633 | 73.9/- | 73.4/- | - | - | -/- | -/- |
| | Conv-BERT | 0.901 | 0.698 | -/80.8 | -/80.7 | 0.593 | 0.700 | -/82.5 | -/82.1 |
| MulT [5]$^{\ddagger}$ | Conv | 0.871 | 0.698 | -/83.0 | -/82.8 | 0.580 | 0.703 | -/82.5 | -/82.3 |
| | Conv-BERT | 0.861 | 0.711 | 81.5/84.1 | 80.6/83.9 | - | - | -/83.5 | -/82.9 |
| MISA [12]$^{\diamond}$ | Conv-BERT | 0.804 | 0.764 | 80.8/82.1 | 80.8/82.0 | 0.568 | 0.724 | 82.6/84.2 | 82.7/84.0 |
| MAG [10]$^{\diamond}$ | Conv-BERT | 0.731 | 0.789 | 82.5/84.3 | 82.6/84.3 | 0.539 | 0.753 | **83.8**/85.2 | **83.7**/85.1 |
| Self-MM [13]$^{\diamond}$ | Conv-BERT | **0.713** | **0.798** | **84.0/86.0** | **84.4/86.0** | 0.530 | 0.765 | 82.8/85.2 | 82.5/85.3 |
| UEGD | Conv | 0.953 | 0.663 | 76.3/77.4 | 76.3/77.4 | 0.598 | 0.683 | 78.9/81.3 | 79.2/81.0 |
| | Conv-BERT | 0.886 | 0.691 | 78.6/79.9 | 78.5/79.9 | 0.543 | 0.748 | 81.2/84.6 | 81.7/84.5 |
| | Enc. output | 0.850 | 0.715 | 79.4/80.8 | 79.3/80.8 | 0.519 | 0.776 | 82.5/85.8 | 82.8/85.7 |
| | Enc. mid-best | 0.828 | 0.748 | 82.0/83.9 | 82.1/84.0 | **0.506** | **0.790** | 82.4/**86.1** | 82.7/**86.0** |
| | Enc. weighted | 0.818 | 0.749 | 80.4/82.3 | 80.4/82.3 | 0.510 | 0.785 | 82.3/85.7 | 82.7/85.6 |

**Table 3**: Performances with single modalities using the UEGD model. Bold means the best performances in each modality.

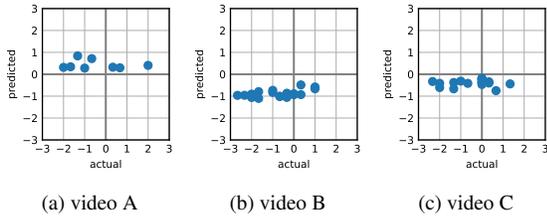| Modality | Feature | CMU-MOSI | | | | CMU-MOSEI | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | MAE↓ | Corr↑ | Acc-2↑ | F-Score↑ | MAE↓ | Corr↑ | Acc-2↑ | F-Score↑ |
| Visual | Facet | **1.431** | **0.147** | **54.0/52.9** | **52.0/51.1** | 0.802 | 0.278 | 68.3/66.1 | 66.1/62.4 |
| | Enc. output | 1.468 | 0.033 | 47.3/45.5 | 41.1/39.5 | 0.771 | 0.403 | 69.6/70.6 | 69.5/69.4 |
| | Enc. mid-best | 1.650 | -0.091 | 44.9/43.4 | 41.5/40.1 | 0.766 | 0.424 | 70.4/**72.2** | **71.0/71.8** |
| | Enc. weighted | 1.630 | 0.074 | 51.2/50.2 | 49.7/48.8 | **0.762** | **0.435** | **70.7**/71.9 | **71.0**/71.3 |
| Acoustic | COVEREP | 1.381 | 0.233 | 55.6/54.7 | 54.5/53.8 | 0.829 | 0.148 | 70.1/63.4 | 61.1/52.1 |
| | Enc. output | 1.326 | 0.326 | 63.4/62.9 | 63.4/63.0 | 0.664 | 0.587 | 74.9/76.1 | 75.0/75.5 |
| | Enc. mid-best | **1.098** | **0.521** | **70.0/70.9** | **70.0/71.0** | **0.605** | **0.666** | **77.4/80.4** | **77.9/80.3** |
| | Enc. weighted | 1.247 | 0.406 | 65.9/66.0 | 65.8/66.0 | 0.635 | 0.629 | 77.0/78.7 | 77.1/78.1 |
| Linguistic | GloVe | 0.963 | 0.634 | 76.6/77.7 | 76.6/77.7 | 0.616 | 0.661 | 78.4/80.1 | 78.6/79.6 |
| | BERT-Base | 0.913 | 0.679 | 78.3/79.6 | 78.2/79.6 | 0.551 | 0.742 | 81.5/84.6 | 81.9/84.5 |
| | Enc. output | 0.854 | 0.715 | 80.6/82.1 | 80.6/82.1 | 0.544 | 0.751 | 81.5/84.9 | 81.9/84.9 |
| | Enc. mid-best | **0.821** | **0.746** | **81.8/83.6** | **81.8/83.6** | 0.540 | **0.760** | **81.7/85.4** | **82.2/85.3** |
| | Enc. weighted | 0.837 | 0.723 | 80.0/81.9 | 80.0/81.9 | **0.539** | 0.758 | 81.1/84.9 | 81.6/84.9 |

### 4.3.2. Evaluations of the Unimodal Model

The prediction results of individual modalities were also compared. The results are shown in Table 3. Except for visual modality in CMU-MOSI, all the encoder-based features outperformed the conventional features in MAE and corr, especially acoustic modality in CMU-MOSEI. With regard to the three encoder-based features, the output of the intermediate layers offers the best performance. One possible reason for the lower performances of the weighted sum is the limited training data. The amounts of training data for the downstream tasks in this work are smaller than those in the previous work [17], e.g., speech and speaker recognition, which may lead to overfitting of the weights of the intermediate outputs.

We conducted further analyses of the encoder output properties and found that the encoder outputs for visual
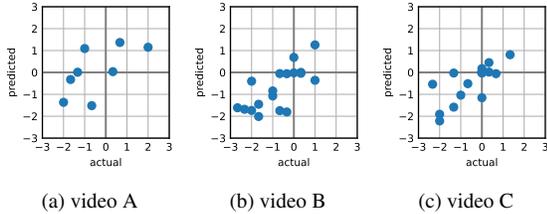


**Fig. 3**: Distributions of the gate weights of the UEGD models on the CMU-MOSEI dataset.

modality may lead to extraction of speaker information rather than sentiment. The prediction examples of the clips in the same videos, i.e., in the same speaker, by the pre-trained encoders in visual and acoustic modalities are shown in Fig. 4 and 5, respectively. As shown in Fig. 4, the model with a

**Fig. 4**: Prediction examples to the clips in the same video (speaker) by the pre-trained visual encoder.
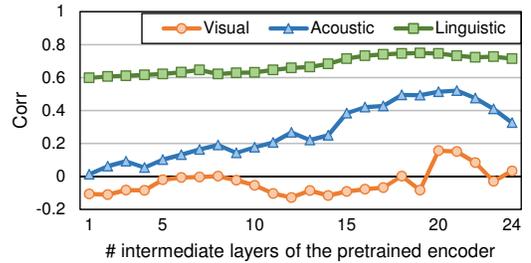


**Fig. 5**: Prediction examples from the clips in the same video (speaker) by the pre-trained acoustic encoder.

pre-trained visual encoder yielded similar prediction results for the same speaker, unlike the pre-trained acoustic encoder. These characteristics appeared in the objective evaluations, total- and intra-video variances of the prediction results for each modality. The results in Table 4 show that the visual encoder yielded much smaller intra-video variances than either the acoustic or linguistic encoder. However, similar properties were observed in the conventional visual features. Further investigation including evaluation of other datasets is required to elucidate the properties of the visual pre-trained encoder.
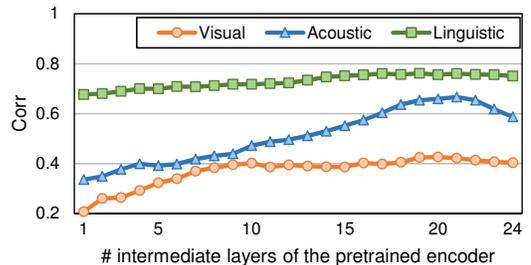
Finally, we discuss the characteristics of the intermediate layer outputs of the pre-trained encoders. The correlation coefficients for each intermediate layer in the two datasets are shown in Fig. 6. The difference in performance of each intermediate layer was small for linguistic modality but significant for the other two modalities. This may be because sentiment information appears as word (low-level) features in linguistic modality, while it appears as action units or speaking styles which are high-level features in visual and acoustic modalities. For the acoustic modality, the highest accuracy was achieved when using around the 20th layer in both datasets. It is considered that it is clearly better to use the middle layer of the second half of the pre-trained encoder than the output of the acoustic modality. For the visual modality, the performances were almost flat after the 9 th layer of the encoder on CMU-MOSEI. This result suggests that there is no intermediate layer strongly associated with sentiment, at least for CLIP ViT-L. Based on these results, best performance is likely to be achieved with the outputs of the second half of the domain-specific pre-trained encoders, especially for the acoustic modality.

**Table 4**: Total and intra-video variances of the predicted values. The conventional linguistic feature was BERT-Base.

| | Modality | CMU-MOSI | | CMU-MOSEI | |
|---|---|---|---|---|---|
| | | Total | Intra | Total | Intra |
| Conv. | Visual | 0.661 | 0.430 | 0.060 | 0.014 |
| | Acoustic | 0.351 | 0.239 | 0.006 | 0.004 |
| | Linguistic | 2.009 | 1.606 | 0.737 | 0.427 |
| Enc. | Visual | 0.984 | 0.353 | 0.230 | 0.016 |
| | Acoustic | 1.412 | 1.178 | 0.677 | 0.361 |
| | Linguistic | 2.005 | 1.612 | 0.750 | 0.414 |
| *Ground Truth* | | 2.523 | 1.775 | 1.229 | 0.606 |



(a) CMU-MOSI



(b) CMU-MOSEI

**Fig. 6**: Performances of the individual intermediate layers of the pre-trained encoders.

## 5. CONCLUSION

This paper investigated the effectiveness and implementation of domain-specific large-scale pre-trained encoders for MSA. The regression model called UEGD was employed to compare the features in unimodal and multimodal scenarios. Three types of features from large-scale pre-trained encoders were compared. The findings from the experiments are as follows. First, the large-scale pre-trained encoder yielded improved sentiment analysis performance in both unimodal and multimodal prediction models when a large amount of labeled training data was available. Second, the pre-trained encoder is particularly effective for acoustic modality. Third, the second half of any of the pre-trained encoders was most effective for MSA, rather than the output. Future work includes further comparisons with other pre-trained encoders on other MSA datasets and the development of a multimodal decoder that can effectively utilize the pre-trained encoders.

# 6. REFERENCES

[1] I. Michael Revina and W. R. Sam Emmanuel, "A survey on human face expression recognition techniques," *Journal of King Saud University - Computer and Information Sciences*, vol. 33, no. 6, pp. 619–628, 2021.

[2] Siddique Latif, Rajib Rana, Sara Khalifa, Raja Jurdak, Junaid Qadir, and Bjoern W. Schuller, "Survey of deep representation learning for speech emotion recognition," *IEEE Trans. on Affective Computing (Early Access)*, 2021.

[3] Marouane Birjali, Mohammed Kasri, and Abderrahim Beni-Hssane, "A comprehensive survey on sentiment analysis: Approaches, challenges and trends," *Knowledge-Based Systems*, vol. 226, pp. 107–134, 2021.

[4] Yansen Wang, Ying Shen, Zhun Liu, Paul Pu Liang, Amir Zadeh, and Louis-Philippe Morency, "Words can shift: Dynamically adjusting word representations using nonverbal behaviors," in *Proc. of AAAI*, 2019, pp. 7216–7223.

[5] Yao-Hung Hubert Tsai, Shaojie Bai, Paul Pu Liang, J. Zico Kolter, Louis-Philippe Morency, and Ruslan Salakhutdinov, "Multimodal transformer for unaligned multimodal language sequences," in *Proc. of ACL*, 2019, pp. 6558–6569.

[6] Amir Zadeh, Minghai Chen, Soujanya Poria, Erik Cambria, and Louis-Philippe Morency, "Tensor fusion network for multimodal sentiment analysis," in *Proc. of EMNLP*, 2017, pp. 1103–1114.

[7] Paul Liang, Ziyin Liu, AmirAli Zadeh, and Louis-Philippe Morency, "Multimodal language analysis with recurrent multistage fusion," in *Proc. of EMNLP*, 2018, pp. 150–161.

[8] Yao-Hung Hubert Tsai, Paul Pu Liang, Amir Zadeh, Louis-Philippe Morency, and Ruslan Salakhutdinov, "Learning factorized multimodal representations," in *Proc. of ICLR*, 2019.

[9] Hai Pham, Paul Pu Liang, Thomas Manzini, Louis-Philippe Morency, and Barnabás Póczos, "Found in translation: Learning robust joint representations by cyclic translations between modalities," in *Proc. of AAAI*, 2019, pp. 6892–6899.

[10] Wasifur Rahman, Md Kamrul Hasan, Sangwu Lee, AmirAli Bagher Zadeh, Chengfeng Mao, Louis-Philippe Morency, and Ehsan Hoque, "Integrating multimodal information in large pretrained transformers," in *Proc. of ACL*, 2020, pp. 2359–2369.

[11] Xianbing Zhao, Yixin Chen, Wanting Li, Lei Gao, and Buzhou Tang, "MAG+: An extended multimodal adaptation gate for multimodal sentiment analysis," in *Proc. of ICASSP*, 2022, pp. 4753–4757.

[12] Devamanyu Hazarika, Roger Zimmermann, and Soujanya Poria, "MISA: Modality-invariant and -specific representations for multimodal sentiment analysis," in *Proc. of ACM MM*, 2020, pp. 1122–1131.

[13] Wenmeng Yu, Hua Xu, Ziqi Yuan, and Jiele Wu, "Learning modality-specific representations with self-supervised multi-task learning for multimodal sentiment analysis," in *Proc. of AAAI*, 2021, pp. 10790–10797.

[14] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei, "Language models are few-shot learners," in *Advances in NeurIPS*, 2020, vol. 33, pp. 1877–1901.

[15] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," in *Advances in NeurIPS*, 2020, vol. 33, pp. 12449–12460.

[16] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever, "Learning transferable visual models from natural language supervision," in *Proc. of ICML*, 2021, vol. 139, pp. 8748–8763.

[17] Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki Kanda, Takuya Yoshioka, Xiong Xiao, Jian Wu, Long Zhou, Shuo Ren, Yanmin Qian, Yao Qian, Jian Wu, Michael Zeng, Xiangzhan Yu, and Furu Wei, "WavLM: Large-scale self-supervised pre-training for full stack speech processing," *IEEE Journal of Selected Topics in Signal Processing*, pp. 1–14, 2022.

[18] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proc. of NAACL*, 2019, pp. 4171–4186.

[19] 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/ACM Trans. on Audio, Speech, and Language Processing*, vol. 29, pp. 3451–3460, 2021.

[20] Ganesh Jawahar, Benoît Sagot, and Djamé Seddah, "What does BERT learn about the structure of language?," in *Proc. of ACL*, 2019, pp. 3651–3657.

[21] Amir Zadeh, Rowan Zellers, Eli Pincus, and Louis-Philippe Morency, "Multimodal sentiment intensity analysis in videos: Facial gestures and verbal messages," *IEEE Intelligent Systems*, vol. 31, no. 16, pp. 82–88, 2016.

[22] AmirAli Bagher Zadeh, Paul Pu Liang, Soujanya Poria, Erik Cambria, and Louis-Philippe Morency, "Multimodal language analysis in the wild: CMU-MOSEI dataset and interpretable dynamic fusion graph," in *Proc. of ACL*, 2018, pp. 2236–2246.

[23] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton, "A simple framework for contrastive learning of visual representations," in *Proc. of ICML*, 2020, pp. 1597–1607.

[24] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, L ukasz Kaiser, and Illia Polosukhin, "Attention is all you need," in *Advances in NeurIPS*, 2017, vol. 30.

[25] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le, "XLNet: Generalized autoregressive pretraining for language understanding," in *Advances in NeurIPS*, 2019, vol. 32.

[26] Kaipeng Zhang, Zhanpeng Zhang, Zhifeng Li, and Yu Qiao, "Joint face detection and alignment using multitask cascaded convolutional networks," *IEEE Signal Processing Letters*, vol. 23, no. 10, pp. 1499–1503, 2016.

[27] Amir Zadeh, Paul Pu Liang, Soujanya Poria, Prateek Vij, Erik Cambria, and Louis-Philippe Morency, "Multi-attention recurrent network for human communication comprehension," in *Proc. of AAAI*, 2018, pp. 5642–5649.

[28] Daniel S. Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D. Cubuk, and Quoc V. Le, "SpecAugment: A Simple Data Augmentation Method for Automatic Speech Recognition," in *Proc. of INTERSPEECH*, 2019, pp. 2613–2617.