

# Hybrid Multimodal Fusion for Humor Detection

Haojie Xu

AHU-IAI AI Joint Laboratory,  
Anhui University  
Institute of Artificial Intelligence,  
Hefei Comprehensive National  
Science Center  
Hefei, China  
e20201139@stu.ahu.edu.cn

Weifeng Liu

AHU-IAI AI Joint Laboratory,  
Anhui University  
Institute of Artificial Intelligence,  
Hefei Comprehensive National  
Science Center  
Hefei, China  
wfeng\_ch@163.com

Jingwei Liu

AHU-IAI AI Joint Laboratory,  
Anhui University  
Institute of Artificial Intelligence,  
Hefei Comprehensive National  
Science Center  
Hefei, China  
ljw578616559@163.com

Mingzheng Li

University of Science and Technology  
of China  
Institute of Artificial Intelligence,  
Hefei Comprehensive National  
Science Center  
Hefei, China  
mingzhengli@mail.ustc.edu.cn

Yu Feng

School of Biomedical Engineering,  
Anhui Medical University  
Institute of Artificial Intelligence,  
Hefei Comprehensive National  
Science Center  
Hefei, China  
fengyu1919@126.com

Yasi Peng

School of Biomedical Engineering,  
Anhui Medical University  
Institute of Artificial Intelligence,  
Hefei Comprehensive National  
Science Center  
Hefei, China  
1558778695@qq.com

Yunwei Shi

AHU-IAI AI Joint Laboratory,  
Anhui University  
Institute of Artificial Intelligence,  
Hefei Comprehensive National  
Science Center  
Hefei, China  
1826719913@qq.com

Xiao Sun\*

Hefei University of Technology  
ZhongJuYuan Intelligent Technology  
Co.,  
Ltd  
Hefei, China  
sunx@iai.ustc.edu.cn

Meng Wang

Hefei University of Technology  
Institute of Artificial Intelligence,  
Hefei Comprehensive National  
Science Center  
Hefei, China  
eric.mengwang@gmail.com

## ABSTRACT

In this paper, we present our solution to the MuSe-Humor sub-challenge of the Multimodal Emotional Challenge (MuSe) 2022. The goal of the MuSe-Humor sub-challenge is to detect humor and calculate AUC from audiovisual recordings of German football Bundesliga press conferences. It is annotated for humor displayed by the coaches. For this sub-challenge, we first build a discriminant model using the transformer module and BiLSTM module, and then propose a hybrid fusion strategy to use the prediction results of each modality to improve the performance of the model. Our experiments demonstrate the effectiveness of our proposed model and hybrid fusion strategy on multimodal fusion, and the AUC of our proposed model on the test set is 0.8972.

\*Corresponding author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

MuSe' 22, October 10, 2022, Lisboa, Portugal

© 2022 Association for Computing Machinery.

ACM ISBN 978-1-4503-9484-0/22/10...\$15.00

<https://doi.org/10.1145/3551876.3554802>

## CCS CONCEPTS

• **Computing methodologies** → **Neural networks**; • **Information systems** → *Multimedia information systems*.

## KEYWORDS

Multimodal Sentiment Analysis; Affective Computing; Humor Detection; Multimodal Fusion

## ACM Reference Format:

Haojie Xu, Weifeng Liu, Jingwei Liu, Mingzheng Li, Yu Feng, Yasi Peng, Yunwei Shi, Xiao Sun, and Meng Wang. 2022. Hybrid Multimodal Fusion for Humor Detection. In *Proceedings of the 3rd International Multimodal Sentiment Analysis Workshop and Challenge (MuSe' 22)*, October 10, 2022, Lisboa, Portugal. ACM, New York, NY, USA, 7 pages. <https://doi.org/10.1145/3551876.3554802>

## 1 INTRODUCTION

Humor is the absurd, unexpected, yet subtle or evocative characteristic of something. It is also a way to help people relieve stress. Not only does it enhance feelings, but it also helps increase overall well-being. As a popular research field in natural language processing (NLP), humor detection has received more attention from scholars [7, 22, 38, 45].

The inconsistency theory of humor argues that humor arises from two or more incongruent but related situations. However, due

to differences in thought, culture and cognition, people's understanding of humor is not the same. This means that humor detection requires a lot of prior knowledge and background information, which brings great challenges for the machine to understand humor. Fortunately, humor often comes from people's interactions. In addition to the text containing the context, one can observe the expression of the speaker in the process of speaking, and we can also obtain the rhythmic clues in his voice. All of this information contributes to the detection of the humorous element. In other words, utilizing the multimodal characteristics of text, acoustic and visual can greatly help the humor detection task [1, 16].

In general, multimodal feature fusion includes early-fusion and late-fusion. For early-fusion, the characteristics of all modalities are fused at first, then they are sent to the model for training. As there are many differences between different modalities, mixing them directly may lead to the underutilization of information. With regard to the late-fusion method, researchers first utilize certain individual models to extract features of each modality. Then, the features being extracted are combined for subsequent tasks. However, the training process may result in the loss of original information. Therefore, we propose a hybrid multimodal fusion model for humor detection, called HMF-MD, which absorbs the advantages of the two fusion methods. First, several discriminant modules are applied to every single modality. Then we combine these features with the original data and send them into another discriminant module for humor detection. This can not only extract the key information inside each modality but also prevent the loss of information effectively. Experiments show that our model achieves good results on the MuSe-Humor sub-challenge task.

Specifically, the main contributions of our work are as follows:

- We propose a discriminative module, which can extract contextual and key information from a single modality.
- We propose a new hybrid fusion strategy that can improve model performance.
- Experiments demonstrate the effectiveness of our proposed model and hybrid fusion strategy on multimodal fusion, and the AUC of our proposed model on the test set is 0.8972.

## 2 RELATED WORKS

### 2.1 Humor Detection

Humor detection has been one of the active areas in the field of affective computing. Humor recognition is to identify whether a sentence or an utterance is humorous or not. There are many ways to be used to identify humor. [38] uses non-neural models to recognize humor. Recently, there are some studies using recurrent neural networks (RNNs) and convolutional neural networks (CNNs) to detect humor[11]. The pre-trained language model has achieved great success in many areas, there are also a lot of works using transformer-based architecture and attention mechanisms to detect humor[4, 26].

Multimodal humor detection is a new area of research in NLP. Multimodal studies from textual, visual and acoustic are the recent research trends. Many works present new Multimodal neural architectures [17, 29, 35], and multimodal fusion approaches [6, 21].

### 2.2 Multimodal Fusion

In a multimodal setting where multiple modalities convey information from different channels, cross-modal fusion is crucial in exploring intra-modality and inter-modality dynamics to mine complementary information. In recent years, multimodal fusion has seen rapid development mainly thanks to the multimodal machine learning community. Earlier multimodal fusion methods fall into two broad categories, i.e. feature-level fusion and decision-level fusion, otherwise known as early and late fusion respectively. Feature-level fusion methods mostly fuse features through concatenation of unimodal representations [27, 30, 36], whereas decision-level fusion methods firstly make a tentative inference for each modality and further fuse them using a voting mechanism [20, 28, 34, 37, 43]. However, these two types of fusion methods cannot effectively explore the inter-modality dynamics. Recently proposed fusion methods can be categorized into several types as follows, i.e., multi-view learning methods [32, 41], word-level fusion methods [15, 35, 42], tensor fusion [24, 40], and hybrid fusion[9, 25]. These fusion techniques are effective in learning inter-modality dynamics compared to feature-level and decision-level fusion and show marked performance gains.

## 3 MULTIMODAL FEATURES

### 3.1 Acoustic Features

All audio files are first normalised to -3 decibels and then converted from stereo to mono, at 16 kHz, 16 bit. Afterwards, we make use of the two well-established machine learning toolkits openSMILE[14] and DeepSpectrum[3] for expert-designed and deep feature extraction from the audio recordings.

**eGeMAPS Feature:** We use the extended Geneva Minimalistic Acoustic Parameter Set (eGeMAPS) feature provided by the organizers of MuSe 2022, which contains 23 acoustic low-level descriptors (LLDs)[13]. The freely and publicly available openSMILE toolkit can be used to extract the eGeMAPS feature. Several statistical functions of the openSMILE toolkit can be directly applied to extract segment-level features with an 88-dimensional vector.

**DeepSpectrum Feature:** The principle of DeepSpectrum is to utilise pre-trained image Convolutional Neural Networks (CNNs) for the extraction of deep features from visual representations (e.g., Mel-spectrograms) of audio signals.

### 3.2 Visual Features

To extract specific image descriptors related to facial expressions, we make use of two CNN architectures: Multi-task Cascaded Convolutional Networks (MTCNN)[44] and VGGface 2[8]. We also provide a set of Facial Action Units (FAUs) obtained from faces of individuals in the datasets.

**MTCNN:** The MTCNN model, pre-trained on the datasets WIDER FACE[39] and CelebA[23], is used to detect faces in the videos. The extracted faces then serve as inputs of the feature extractors VGGface 2.

**VGGface 2:** The purpose of VGGface 2 is to compute general facial features for the previously extracted faces. VGGface 2 is a dataset for the task of face recognition. It contains 3.3 million faces of about 9,000 different persons. We use a ResNet50[18] trained

on VGGface 2 and detach its classification layer, resulting in a 512-dimensional feature vector output referred to as VGGface 2.

**FAUs:** Facial Action Units (FAUs) denote the presence of facial muscle movements that are commonly used for describing and classifying expressions. It can be extracted by OpenFace toolkit [5]. We only use the FAU intensity features provided by the organizers.

### 3.3 Textual Features

**Bert:** As the organizer did[2, 10], We employ a German version of the BERT (Bidirectional Encoder Representations from Transformers (BERT)[12]) model. No further fine-tuning is applied. For Passau Spontaneous Football Coach Humor (Passau-SFCH), the dataset used for the humor detection sub-challenge, we extract the BERT token embeddings. Additionally, we obtain 768-dimensional sentence embeddings for all texts in Passau-SFCH by using the encodings of BERT's token. In all cases, we average the embeddings provided by the last 4 layers of the BERT model, following[31].

## 4 MUSE-HUMOR METHOD

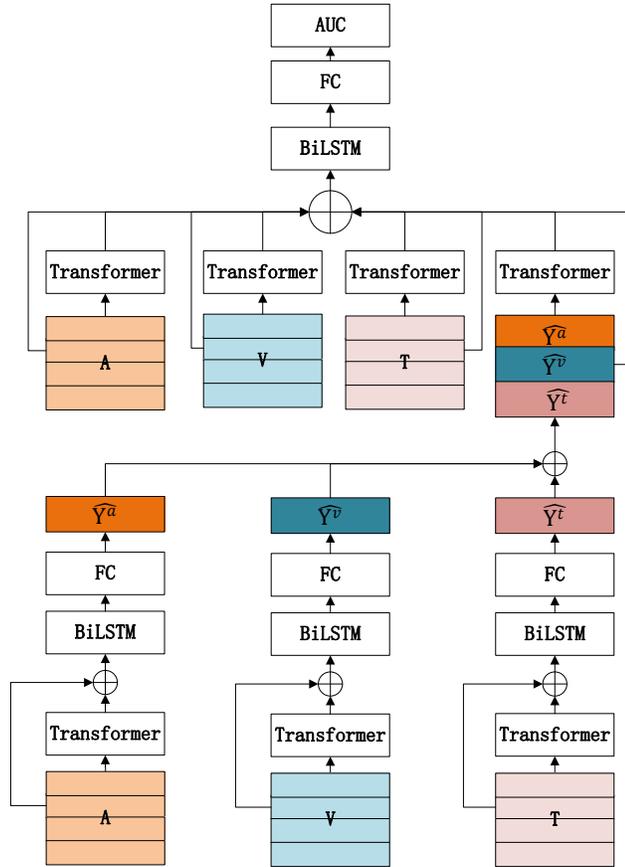


Figure 1: Overview of the architecture used in the MuSe-Humor sub-challenge.

In this section, we will introduce the main components of our approach, and its overall framework is shown in Figure 1. We propose a hybrid multimodal fusion model for humor detection, called HMF-MD. The training process of HMF-MD consists of two stages: (1) Unimodal Discrimination Stage, which obtains the optimal discriminant distribution of the input sequence corresponding to each modality through our discriminant module. (2) Hybrid Multimodal Fusion Stage, which utilizes the Hybrid Multimodal Fusion strategy to fuse the initial input of various modalities and the discriminant output of every single modality in the previous stage, and finally performs humor detection through our discriminant module.

### 4.1 Unimodal Discrimination Stage

For this humor detection sub-challenge, the data is not monomodal, so interactions of multimodal information are inevitably considered. However, for data of a particular modality, there are internal correlations specific to that modality, so utilizing single modality data for the effective fusion of multimodal information is essential.

To capture this correlation within a certain modality, we propose the discriminant module shown in the bottom half of Figure 1. Given an input sequence  $X^m = \{x_1^m, x_2^m, \dots, x_L^m\}$  for a specific modality  $m \in \{A, V, T\}$ , where  $L$  is the sequence length, we first feed it into a Transformer layer[33] to capture the interaction between each element in the sequence and other elements and retain more information about the element itself through a residual connection.

$$H^m = \mathcal{T}(X^m) \oplus X^m \quad (1)$$

where  $H^m = \{h_1^m, h_2^m, \dots, h_L^m\}$  is the hidden representation sequence output by Eq.1,  $\mathcal{T}(\cdot)$  denotes the calculation process in the Transformer layer[33], and  $m \in \{A, V, T\}$  represents a particular modality.

Since the data of the humor detection sub-challenge has an obvious contextual relationship, we then send the obtained hidden representation sequence into a bidirectional LSTM [19] to capture this contextual information and take the hidden output of the last element in the sequence as the hidden representation of this sequence.

$$\tilde{H}^m = \overrightarrow{\mathcal{B}}(H^m)[-1] \parallel \overleftarrow{\mathcal{B}}(H^m)[0] \quad (2)$$

where  $\tilde{H}^m$  is output of the BiLSTM layer,  $\overrightarrow{\mathcal{B}}(\cdot)$  and  $\overleftarrow{\mathcal{B}}(\cdot)$  denote the forward and backward calculation process of LSTM,  $[-1]$  and  $[0]$  represent the last and first elements of the corresponding output sequence, and  $\parallel$  denotes the concatenate operation.

Finally, we send it to the last component of the discrimination module, a fully connected layer, and its output is the representation that captures the internal correlation of the corresponding modality, i.e. the discriminant output of each modality.

$$\hat{Y}^m = \sigma(W^m \tilde{H}^m + b^m) \quad (3)$$

where  $W^m$  and  $b^m$  are the modality-specific weight matrix and the bias term, respectively, and  $\sigma$  is the Sigmoid activation function.

For the input sequence of each modality, we carry out a complete training process to obtain the optimal discriminant output.

### 4.2 Hybrid Multimodal Fusion Stage

In order to integrate the information of various modality more effectively and improve the advantages brought by information

**Table 1: Statistics information of the Passau-SFCH dataset.**

Partition	Coaches	Labels	Duration
Train	4	14025	3 :52 :44
Development	3	11320	3 :08 :12
Test	3	14143	3 :55 :41
Sum	10	39488	10 :56 :37

complementarity between modalities, we propose the hybrid fusion strategy at this stage.

Firstly, we concatenate the discriminant outputs of each modality in the first stage and take them as the input of the second stage.

$$X^{\hat{Y}} = \parallel_{m \in \{A, V, T\}} \hat{Y}^m \quad (4)$$

However, as described in Section 4.1, these discriminant outputs are modality-specific internal correlations, and simply using them for the final detection task cannot take full advantage of the complementary effects between the various modalities. Therefore, in addition to the output of the first stage, the original input of each modality is also taken as the input feature of the second stage, so that the model can effectively capture the complementary information of each modality.

The main structure of the second stage is still the discrimination module introduced in Section 4.1. The difference is that in the second stage, each type of input first passes through its own transformer layer first, and then the combined output is sent to the shared subsequent layers for humor detection.

$$Z = \parallel_{m \in \{A, V, T, \hat{Y}\}} \mathcal{T}^m(X^m) \oplus X^m \quad (5)$$

$$\tilde{Z} = \vec{\mathcal{B}}(Z)[-1] \parallel \overleftarrow{\mathcal{B}}(Z)[0] \quad (6)$$

$$\hat{Z} = \sigma(W\tilde{Z} + b) \quad (7)$$

### 4.3 Wrapped BCELoss

In both stage, we train our model using the Wrapped BCELoss, which can be formatted as follows:

$$\ell = \frac{1}{N} \sum_{n=1}^N l_n \quad (8)$$

$$l_n = -w_n [y_n \cdot \log x_n + (1 - y_n) \cdot \log(1 - x_n)] \quad (9)$$

where  $N$  is the total number of samples,  $l_n$  is the corresponding loss of the  $n$ -th sample, and  $x_n$  and  $y_n$  are the ground-truth label and the predicted label for the  $n$ -th sample.

## 5 EXPERIMENTS

### 5.1 Dataset

In MuSe 2022, three datasets are provided for three different sub-challenges. This paper mainly focuses on the study of humor detection sub-challenge. The dataset for the humor detection sub-challenge is the novel Passau Spontaneous Football Coach Humor (Passau-SFCH) database. It comprises audiovisual recordings of German football Bundesliga press conferences. It is annotated for humor displayed by the coaches. For the challenge, binary labeling

**Table 2: The AUC performance on the development set of the MuSe-Humor sub-challenge. ‘M’ denotes the used modality. ‘A’, ‘V’ and ‘T’ represent the audio, video, and text modality. ‘BiLSTM’ means the official baseline.**

Feature	M	Model	AUC
eGeMAPS	A	BiLSTM	0.6861
eGeMAPS	A	Ours	0.6912
DeepSpectrum	A	BiLSTM	0.7149
DeepSpectrum	A	Ours	0.7187
FAU	V	BiLSTM	0.9071
FAU	V	Ours	0.9050
VGGface 2	V	BiLSTM	0.9253
VGGface 2	V	Ours	0.9331
BERT	T	BiLSTM	0.8270
BERT	T	Ours	0.8261

(presence or absence of humor) is provided. Each label in Passau-SFCH dataset is predicted based on all feature vectors belonging to the corresponding 2 s window. The statistics information is shown in Table 1.

### 5.2 Experimental Setup

We implement all of our models with the PyTorch toolkit in the MuSe-Humor sub-challenges. For the MuSe-Humor sub-challenge, the proposed model consists of the transformer layer, a bidirectional LSTM layer, and a fully connected layer. The number of layers used by all transformer layers is 1, and the number of heads used by all Multi-head attention is 1. The number of hidden sizes in the LSTM layer is 32, and the number of the bidirectional LSTM layer is set to 2. The Adam optimizer with a learning rate of 0.001 is applied, and the batch size is set to 32. Same as the baseline, we train the model for a maximum of 100 epochs and stop training if the validation AUC does not increase for 3 consecutive epochs. For the hybrid fusion model, we only use 0.4 for the dropout of the linear layer and 0.2 for the rest of the modules.

### 5.3 Unimodal Results

We first evaluate the performance of each modality we used in the MuSe-Humor sub-challenge. To verify the effectiveness of the proposed model, several experiments are conducted. The experiment results are given in Table 2. From the table 2, we can conclude that: 1) The ‘Ours’ model achieves the best performance or performance close to ‘BiLSTM’ when using audio and video modalities for classification. 2) ‘BiLSTM’ model outperforms ‘Ours’ models when using text for classification. We believe that the reason is that the feature continuity has deteriorated during feature extraction and processing, and better sequence information can be obtained by directly feeding it into BiLSTM. Nevertheless, the performance difference between ‘Ours’ and ‘BiLSTM’ is very small. 3) On the whole, ‘Ours’ can achieve good performance in the independent use of the three modalities.

**Table 3: AUC performance of different modalities using hybrid fusion strategy on the development set in the MuSe-Humor sub-challenge.**  $\hat{Y}$  represents the discriminative output of the corresponding mode.  $\hat{Y}^m$  represents the discriminative output obtained by using BiLSTM for the corresponding modality.  $\hat{Y}^{\hat{M}}$  represents the discriminative output obtained by using our proposed model for the corresponding modality.

Feature	M	Model	AUC
DeepSpectrum + VGGface	A+V	BiLSTM	0.8252
DeepSpectrum + VGGface	A+V	Ours	0.9306
DeepSpectrum + VGGface	A+V+ $\hat{Y}^{\hat{A},V}$	Ours	0.9415
DeepSpectrum + BERT	A+T	BiLSTM	0.8901
DeepSpectrum + BERT	A+T	Ours	0.8303
DeepSpectrum + BERT	A+T+ $\hat{Y}^{\hat{A},T}$	Ours	0.8508
VGGface 2 + BERT	V+T	BiLSTM	0.8908
VGGface 2 + BERT	V+T	Ours	0.9461
VGGface 2 + BERT	V+T + $\hat{Y}^{\hat{V},T}$	Ours	0.9483
DeepSpectrum + VGGface 2 +BERT	A+V+T	BiLSTM	0.9033
DeepSpectrum + VGGface 2 +BERT	A+V+T	Ours	0.9534
DeepSpectrum + VGGface 2 +BERT	A+V+T+ $\hat{Y}^{\hat{A},\hat{V},T}$	Ours	0.9552
DeepSpectrum + VGGface 2 +BERT	A+V+T+ $\hat{Y}^{\hat{A},\hat{V},t}$	Ours	0.9570
DeepSpectrum + VGGface 2 +BERT	A+V+T+ $\hat{Y}^{\hat{A},\hat{V},t}$	Ours	0.9573

## 5.4 Multimodal Results

When performing multi-modal feature fusion, we select the same modal features as the official baseline for fusion. Compared with the official baseline, we also get better results on the single modality for text features (BERT), audio features (DeepSpectrum) and visual features (VGGface 2). Table 3 shows the AUC performance of different modalities with and without the hybrid fusion strategy on the development set of the MuSe-Humor sub-challenge. To better verify the effectiveness of our proposed hybrid fusion strategy, we list the discriminative output  $\hat{Y}$  as a new modality in Table 3. For example, ‘A+V+T’ means that the features of these three modalities are directly sent to different Transformer layers without using a hybrid fusion strategy, and then the connected features are sent to the BiLSTM layer. ‘A+V+T+ $\hat{Y}^{\hat{A},\hat{V},T}$ ’ means using a hybrid fusion strategy, that is, the discriminative output  $\hat{Y}^{\hat{A},\hat{V},T}$  is obtained first, and then the three modalities and the discriminative outputs obtained from the three modalities are sent to different Transformer layers.

When we perform feature fusion of the three modalities, we also utilize BiLSTM to obtain discriminative outputs for textual modalities and acoustic modalities. In Table 3, ‘A+V+T+ $\hat{Y}^{\hat{A},\hat{V},t}$ ’ means that BiLSTM obtains the discriminative output of text modality, and our proposed model also obtains the discriminative output of acoustic and visual modality. ‘A+V+T+ $\hat{Y}^{\hat{A},\hat{V},t}$ ’ means that BiLSTM obtains the discriminative output of textual and acoustic modality, and our proposed model also obtains the discriminative output of visual modality.

From Table 3, we can find that 1) multi-modal features can achieve better performance than single-modal features. 2) In the results of our model, the results with the hybrid fusion strategy are all better than those without the hybrid fusion strategy. 3) When

**Table 4: The best submission results of our proposed method in the MuSe-Humor sub-challenges.**

Model	M	Development	Test
Baseline	V	0.9253	0.8480
Baseline	A+V+T	0.9033	0.7973
Ours	A+V+T	0.9534	0.8559
Ours	A+V+T+ $\hat{Y}^{\hat{A},\hat{V},T}$	0.9552	0.8823
Ours	A+V+T+ $\hat{Y}^{\hat{A},\hat{V},t}$	0.9570	0.8882
Ours	A+V+T+ $\hat{Y}^{\hat{A},\hat{V},t}$	0.9573	0.8945
Ours	VOTE	-	0.8972

multi-modal feature fusion is performed, better performance can be obtained by the fusion of video and text features or audio and video features, which is consistent with the phenomenon that the video feature results are optimal in a single modality. 4) When the audio and text features are fused, we find that the model results are lower than the official baseline BiLSTM, which we suspect is related to feature extraction and processing. Because each label is predicted based on all feature vectors belonging to the corresponding 2s window, the text features and audio features of a whole sentence may be truncated, and the direct use of BiLSTM can better maintain the temporal continuity of features. 5) Using BiLSTM to obtain the discriminative output of text modality and audio modality can improve the model performance to a certain extent.

## 5.5 Submission Results

Table 4 shows the best submission results of the proposed method in the MuSe-Humor sub-challenge. The official baseline achieves the best AUC results on the test set on video modality features, and

the optimal AUC of our proposed model is 0.0343 higher than the official baseline. Compared with the official baseline results using three modalities, our proposed model outperforms the baseline by 0.085 in AUC. Further, we also use ‘A+V+T+Y<sup>A,V,t</sup>’, ‘A+V+T+Y<sup>A,V,t</sup>’, and ‘A+V+T+Y<sup>A,V,T</sup>’s prediction results to vote, and obtain the average of the predicted results. After submitting the average we reached 0.8972 on AUC. In addition, from the table 4, we also found that adding a hybrid fusion strategy can also improve the recognition ability of the model, so that the model can achieve better results on the test set. When we use the hybrid fusion strategy, the model has a small performance improvement on the validation set, but the test set performance is greatly improved. Overall, the AUC results achieve a decent performance improvement over the official baseline. However, the AUC of the test set is still relatively different from the AUC of the validation set. We conjecture that the distribution of test and validation sets may be different and that different coaches have different habits of expressing humor, so the model may have an overfitting problem.

## 6 CONCLUSIONS

In this paper, we present our solutions for the MuSe-Humor sub-challenge of Multi-modal Sentiment Challenge (MuSe) 2022. For the MuSe-Humor sub-challenge, fusing the textual features (BERT) with the audio features (DeepSpectrum) and the visual features (VGGface 2) proved useful for calculating AUC. Also, our proposed hybrid fusion strategy can further improve the model performance.

The proposed method shows promising prospects for future improvements. First, more features of different modalities can be used in this sub-challenge. Then, the next step can be to explore the number of Transformer layers used by different modalities and other methods to increase the superiority and generalizability of the model. Finally, ways to maintain the temporal continuity of features during extraction and processing can be explored.

## REFERENCES

- [1] Malak Abdullah, Mirsad Hadzikadic, and Samira Shaikh. 2018. SEDAT: Sentiment and Emotion Detection in Arabic Text Using CNN-LSTM Deep Learning. In *17th IEEE International Conference on Machine Learning and Applications, ICMLA 2018, Orlando, FL, USA, December 17-20, 2018*, M. Arif Wani, Mehmed M. Kantardzic, Moamar Sayed Mouchaweh, João Gama, and Edwin Lughofer (Eds.). IEEE, 835–840. <https://doi.org/10.1109/ICMLA.2018.00134>
- [2] Shahin Amiriparian, Lukas Christ, Andreas König, Eva-Maria Meßner, Alan Cowen, Erik Cambria, and Björn W. Schuller. 2022. MuSe 2022 Challenge: Multimodal Humour, Emotional Reactions, and Stress. In *Proceedings of the 30th ACM International Conference on Multimedia (MM'22), October 10-14, 2022, Lisbon, Portugal*. Association for Computing Machinery, Lisbon, Portugal. 3 pages, to appear.
- [3] Shahin Amiriparian, Maurice Gerczuk, Sandra Ottl, Nicholas Cummins, Michael Freitag, Sergey Pugachevskiy, Alice Baird, and Björn Schuller. 2017. Snore sound classification using image-based deep spectrum features. (2017).
- [4] Issa Annamoradnejad and Gohar Zoghi. 2020. Colbert: Using bert sentence embedding for humor detection. *arXiv preprint arXiv:2004.12765* (2020).
- [5] Tadas Baltrušaitis, Peter Robinson, and Louis-Philippe Morency. 2016. Openface: an open source facial behavior analysis toolkit. In *2016 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 1–10.
- [6] Elham J Barezi and Pascale Fung. 2018. Modality-based factorization for multi-modal fusion. *arXiv preprint arXiv:1811.12624* (2018).
- [7] Davide Buscaldi and Paolo Rosso. 2007. Some Experiments in Humour Recognition Using the Italian Wikiquote Collection. In *Applications of Fuzzy Sets Theory, 7th International Workshop on Fuzzy Logic and Applications, WILF 2007, Camogli, Italy, July 7-10, 2007, Proceedings (Lecture Notes in Computer Science, Vol. 4578)*, Francesco Masulli, Sushmita Mitra, and Gabriella Pasi (Eds.). Springer, 464–468. [https://doi.org/10.1007/978-3-540-73400-0\\_58](https://doi.org/10.1007/978-3-540-73400-0_58)
- [8] Qiong Cao, Li Shen, Weidi Xie, Omkar M Parkhi, and Andrew Zisserman. 2018. Vggface2: A dataset for recognising faces across pose and age. In *2018 13th IEEE international conference on automatic face & gesture recognition (FG 2018)*. IEEE, 67–74.
- [9] Haifeng Chen, Yifan Deng, Shiwen Cheng, Yixuan Wang, Dongmei Jiang, and Hichem Sahli. 2019. Efficient Spatial Temporal Convolutional Features for Audiovisual Continuous Affect Recognition. In *Proceedings of the 9th International on Audio/Visual Emotion Challenge and Workshop (Nice, France) (AVEC '19)*. Association for Computing Machinery, New York, NY, USA, 19–26. <https://doi.org/10.1145/3347320.3357690>
- [10] Lukas Christ, Shahin Amiriparian, Alice Baird, Panagiotis Tzirakis, Alexander Kathan, Niklas Müller, Lukas Stappen, Eva-Maria Meßner, Andreas König, Alan Cowen, Erik Cambria, and Björn W. Schuller. 2022. The MuSe 2022 Multimodal Sentiment Analysis Challenge: Humor, Emotional Reactions, and Stress. In *Proceedings of the 3rd Multimodal Sentiment Analysis Challenge*. Association for Computing Machinery, Lisbon, Portugal. Workshop held at ACM Multimedia 2022, to appear.
- [11] Luke De Oliveira and Alfredo L Rodrigo. 2015. Humor detection in yelp reviews. Retrieved on December 15 (2015), 2019.
- [12] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [13] Florian Eyben, Klaus R Scherer, Björn W Schuller, Johan Sundberg, Elisabeth André, Carlos Busso, Laurence Y Devillers, Julien Epps, Petri Laukka, Shrikanth S Narayanan, et al. 2015. The Geneva minimalist acoustic parameter set (GeMAPS) for voice research and affective computing. *IEEE transactions on affective computing* 7, 2 (2015), 190–202.
- [14] Florian Eyben, Martin Wöllmer, and Björn Schuller. 2010. Opensmile: the munich versatile and fast open-source audio feature extractor. In *Proceedings of the 18th ACM international conference on Multimedia*. 1459–1462.
- [15] Yue Gu, Xinyu Li, Kaixiang Huang, Shiyu Fu, Kangning Yang, Shuhong Chen, Moliang Zhou, and Ivan Marsic. 2018. Human Conversation Analysis Using Attentive Multimodal Networks with Hierarchical Encoder-Decoder. In *2018 ACM Multimedia Conference on Multimedia Conference, MM 2018, Seoul, Republic of Korea, October 22-26, 2018*, Susanne Boll, Kyoung Mu Lee, Jiebo Luo, Wenwu Zhu, Hyeran Byun, Chang Wen Chen, Rainer Lienhart, and Tao Mei (Eds.). ACM, 537–545. <https://doi.org/10.1145/3240508.3240714>
- [16] Md. Kamrul Hasan, Wasifur Rahman, AmirAli Bagher Zadeh, Jianyuan Zhong, Md. Iftekhar Tanveer, Louis-Philippe Morency, and Mohammed (Ehsan) Hoque. 2019. UR-FUNNY: A Multimodal Language Dataset for Understanding Humor. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojuan Wan (Eds.). Association for Computational Linguistics, 2046–2056. <https://doi.org/10.18653/v1/D19-1211>
- [17] Devamanyu Hazarika, Soujanya Poria, Amir Zadeh, Erik Cambria, Louis-Philippe Morency, and Roger Zimmermann. 2018. Conversational memory network for emotion recognition in dyadic dialogue videos. In *Proceedings of the conference. Association for Computational Linguistics. North American Chapter. Meeting, Vol. 2018*. NIH Public Access, 2122.
- [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [19] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long Short-Term Memory. *Neural Comput.* 9, 8 (1997), 1735–1780. <https://doi.org/10.1162/neco.1997.9.8.1735>
- [20] Onno Kampman, Elham J. Barezi, Dario Bertero, and Pascale Fung. 2018. Investigating Audio, Visual, and Text Fusion Methods for End-to-End Automatic Personality Prediction. *CoRR abs/1805.00705* (2018). [arXiv:1805.00705](http://arxiv.org/abs/1805.00705)
- [21] Paul Pu Liang, Ziyin Liu, Amir Zadeh, and Louis-Philippe Morency. 2018. Multimodal language analysis with recurrent multistage fusion. *arXiv preprint arXiv:1808.03920* (2018).
- [22] Lizhen Liu, Donghai Zhang, and Wei Song. 2018. Modeling Sentiment Association in Discourse for Humor Recognition. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 2: Short Papers*, Iryna Gurevych and Yusuke Miyao (Eds.). Association for Computational Linguistics, 586–591. <https://doi.org/10.18653/v1/P18-2093>
- [23] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. 2015. Deep learning face attributes in the wild. In *Proceedings of the IEEE international conference on computer vision*. 3730–3738.
- [24] Zhun Liu, Ying Shen, Varun Bharadhwaj Lakshminarasimhan, Paul Pu Liang, Amir Zadeh, and Louis-Philippe Morency. 2018. Efficient Low-rank Multimodal Fusion With Modality-Specific Factors. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 1: Long Papers*, Iryna Gurevych and Yusuke Miyao (Eds.). Association for Computational Linguistics, 2247–2256. <https://doi.org/10.18653/v1/P18-1209>

- [25] Ziyu Ma, Fuyan Ma, Bin Sun, and Shutao Li. 2021. Hybrid Multimodal Fusion for Dimensional Emotion Recognition. In *Proceedings of the 2nd on Multimodal Sentiment Analysis Challenge (Virtual Event, China) (MuSe '21)*. Association for Computing Machinery, New York, NY, USA, 29–36. <https://doi.org/10.1145/3475957.3484457>
- [26] Jihang Mao and Wanli Liu. 2019. A BERT-based Approach for Automatic Humor Detection and Scoring. In *IberLEF@SEPLN*. 197–202.
- [27] 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, ICMI 2011, Alicante, Spain, November 14-18, 2011*, Hervé Bourlard, Thomas S. Huang, Enrique Vidal, Daniel Gatica-Perez, Louis-Philippe Morency, and Nicu Sebe (Eds.). ACM, 169–176. <https://doi.org/10.1145/2070481.2070509>
- [28] Behnaz Nojavanasghari, Deepak Gopinath, Jayanth Koushik, Tadas Baltrusaitis, and Louis-Philippe Morency. 2016. Deep multimodal fusion for persuasiveness prediction. In *Proceedings of the 18th ACM International Conference on Multimodal Interaction, ICMI 2016, Tokyo, Japan, November 12-16, 2016*, Yukiko I. Nakano, Elisabeth André, Toyooki Nishida, Louis-Philippe Morency, Carlos Busso, and Catherine Pelachaud (Eds.). ACM, 284–288. <https://doi.org/10.1145/2993148.2993176>
- [29] Soujanya Poria, Erik Cambria, Devamanyu Hazarika, Navonil Mazumder, Amir Zadeh, and Louis-Philippe Morency. 2017. Multi-level multiple attentions for contextual multimodal sentiment analysis. In *2017 IEEE International Conference on Data Mining (ICDM)*. IEEE, 1033–1038.
- [30] Viktor Rozgic, Sankaranarayanan Ananthakrishnan, Shirin Saleem, Rohit Kumar, and Rohit Prasad. 2012. Ensemble of SVM trees for multimodal emotion recognition. In *Asia-Pacific Signal and Information Processing Association Annual Summit and Conference, APSIPA 2012, Hollywood, CA, USA, December 3-6, 2012*. IEEE, 1–4. <https://ieeexplore.ieee.org/document/6411794/>
- [31] Licai Sun, Zheng Lian, Jianhua Tao, Bin Liu, and Mingyue Niu. 2020. Multi-modal continuous dimensional emotion recognition using recurrent neural network and self-attention mechanism. In *Proceedings of the 1st International on Multimodal Sentiment Analysis in Real-life Media Challenge and Workshop*. 27–34.
- [32] Shiliang Sun. 2013. A survey of multi-view machine learning. *Neural Comput. Appl.* 23, 7-8 (2013), 2031–2038. <https://doi.org/10.1007/s00521-013-1362-6>
- [33] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (Eds.). 5998–6008. <https://proceedings.neurips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html>
- [34] Haoan 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)*. 949–954. <https://doi.org/10.1109/ICME.2017.8019301>
- [35] 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 *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019*. AAAI Press, 7216–7223. <https://doi.org/10.1609/aaai.v33i01.33017216>
- [36] Martin Wöllmer, Felix Weninger, Tobias Knaup, Björn W. Schuller, Congkai Sun, Kenji Sagae, and Louis-Philippe Morency. 2013. YouTube Movie Reviews: Sentiment Analysis in an Audio-Visual Context. *IEEE Intell. Syst.* 28, 3 (2013), 46–53. <https://doi.org/10.1109/MIS.2013.34>
- [37] Chung-Hsien Wu and Wei-Bin Liang. 2011. Emotion Recognition of Affective Speech Based on Multiple Classifiers Using Acoustic-Prosodic Information and Semantic Labels. *IEEE Trans. Affect. Comput.* 2, 1 (2011), 10–21. <https://doi.org/10.1109/T-AFFC.2010.16>
- [38] Diyi Yang, Alon Lavie, Chris Dyer, and Eduard H. Hovy. 2015. Humor Recognition and Humor Anchor Extraction. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, EMNLP 2015, Lisbon, Portugal, September 17-21, 2015*, Lluís Màrquez, Chris Callison-Burch, Jian Su, Daniele Pighin, and Yuval Marton (Eds.). The Association for Computational Linguistics, 2367–2376. <https://doi.org/10.18653/v1/d15-1284>
- [39] Shuo Yang, Ping Luo, Chen-Change Loy, and Xiaoou Tang. 2016. Wider face: A face detection benchmark. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 5525–5533.
- [40] Amir Zadeh, Minghai Chen, Soujanya Poria, Erik Cambria, and Louis-Philippe Morency. 2017. Tensor Fusion Network for Multimodal Sentiment Analysis. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP 2017, Copenhagen, Denmark, September 9-11, 2017*, Martha Palmer, Rebecca Hwa, and Sebastian Riedel (Eds.). Association for Computational Linguistics, 1103–1114. <https://doi.org/10.18653/v1/d17-1115>
- [41] Amir Zadeh, Paul Pu Liang, Navonil Mazumder, Soujanya Poria, Erik Cambria, and Louis-Philippe Morency. 2018. Memory Fusion Network for Multi-view Sequential Learning. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence (AAAI-18), the 30th Innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*, Sheila A. McIlraith and Kilian Q. Weinberger (Eds.). AAAI Press, 5634–5641. <https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/17341>
- [42] Amir Zadeh, Paul Pu Liang, Soujanya Poria, Prateek Vij, Erik Cambria, and Louis-Philippe Morency. 2018. Multi-attention Recurrent Network for Human Communication Comprehension. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence (AAAI-18), the 30th Innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*, Sheila A. McIlraith and Kilian Q. Weinberger (Eds.). AAAI Press, 5642–5649. <https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/17390>
- [43] Amir Zadeh, Rowan Zellers, Eli Pincus, and Louis-Philippe Morency. 2016. MOSI: Multimodal Corpus of Sentiment Intensity and Subjectivity Analysis in Online Opinion Videos. CoRR abs/1606.06259 (2016). arXiv:1606.06259 <http://arxiv.org/abs/1606.06259>
- [44] 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 (2016), 1499–1503.
- [45] Dingyu Zhu. 2019. Humor robot and humor generation method based on big data search through IOT. *Clust. Comput.* 22, Supplement (2019), 9169–9175. <https://doi.org/10.1007/s10586-018-2097-z>