

Rethinking the constraints of multimodal fusion: case study in Weakly-Supervised Audio-Visual Video Parsing

Jianning Wu^{1,b} Zhuqing Jiang^{1,2,a,b} Shiping Wen³ Aidong Men¹ Haiying Wang¹

^a Corresponding author

^b The first two authors contribute equally to this work

¹ School of Artificial Intelligence, Beijing University of Posts and Telecommunications, Beijing, China

² Beijing Key Laboratory of Network System and Network Culture, Beijing University of Posts and Telecommunications, Beijing, China

³ Australian AI Institute, Faculty of Engineering and Information Technology, University of Technology Sydney, Australia

{jianningwu, jiangzhuqing}@bupt.edu.cn

Abstract

For multimodal tasks, a good feature extraction network should extract information as much as possible and ensure that the extracted feature embedding and other modal feature embedding have an excellent mutual understanding. The latter is often more critical in feature fusion than the former. Therefore, selecting the optimal feature extraction network collocation is a very important subproblem in multimodal tasks. Most of the existing studies ignore this problem or adopt an ergodic approach. This problem is modeled as an optimization problem in this paper. A novel method is proposed to convert the optimization problem into an issue of comparative upper bounds by referring to the general practice of extreme value conversion in mathematics. Compared with the traditional method, it reduces the time cost.

Meanwhile, aiming at the common problem that the feature similarity and the feature semantic similarity are not aligned in the multimodal time-series problem, we refer to the idea of contrast learning and propose a multimodal time-series contrastive loss (MTSC).

Based on the above issues, We demonstrated the feasibility of our approach in the audio-visual video parsing task. Substantial analyses verify that our methods promote the fusion of different modal features.

1. Introduction

Audio-Visual Video Parsing (AVVP) [33] has a wide potential application in downstream video understanding tasks (such as monitoring analysis, video summarization, and retrieval). It is a newly introduced multi-modal task that involves detecting and localizing occurrences of events within the audio and visual streams of a video. It directly contributes to audio-visual source separation, espe-

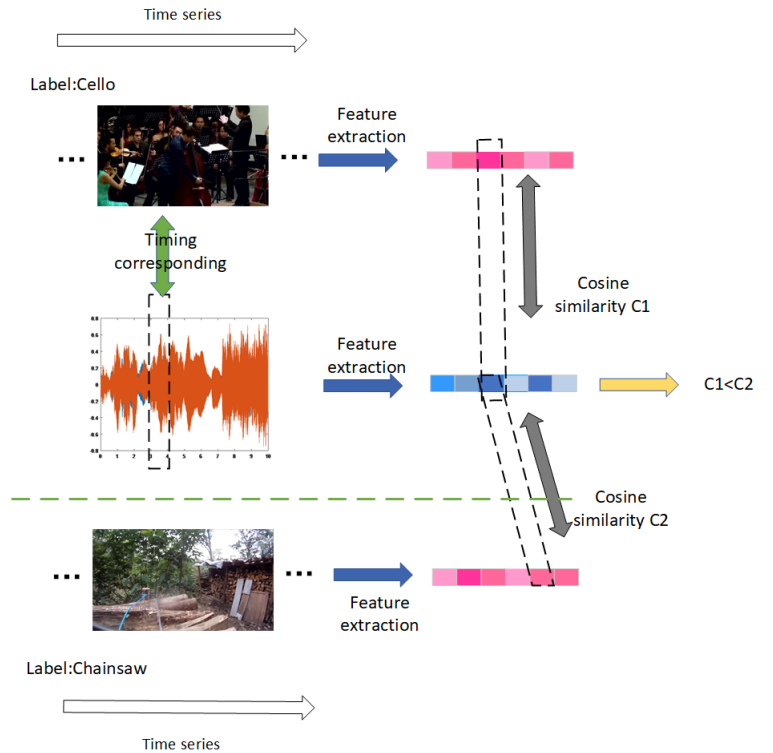


Figure 1. $C2$ is the maximum similarity of the audio fragments in the figure. $C1$ is the similarity between the image frame of the corresponding segment and the audio at the same time; $C2$ is the similarity between the image frames and the audio of different video clips.

cially when the sources of audio are occluded in the video.

In general, the standard process of the multimodal problem [15, 18, 28, 30, 39] could be divided into three steps: feature extraction/embedding, feature fusion, and subsequent feature processing. Among them, the most

characteristic and challenging step for multimodal missions is multimodal fusion[5].(As shown in Figure 2) Obtained through the first two parts, the fusion feature plays an important role in the working performance of subsequent feature processing.

However, when designing network structure and loss function, the existing AVVP method[33] only takes the characteristics of TAL task into consideration and pays less attention to the difficulty and complexity of multimodal feature fusion. Considering it, we propose that the main AVVP methods have two problems:

- 1.Lack of a practical approach to select pre-trained feature extraction network collocation effectively.
- 2.Ignoring the difference between feature similarity and feature semantic similarity.

The first problem exists not only in only AVVP but also in most of the multimodal tasks as long as the extraction network is needed. The inevitability of this problem is mainly based on the following fact: 1. Except for the difference inside the raw multi-mode data, the feature extraction network also significantly impacts subsequent feature fusion, which is ignored and underestimated in common sense; 2. there are plenty of pre-trained feature extraction networks but few methods for selecting them, except traverse; 3. Meanwhile, in today’s most popular multimodal presentation learning task, the cost of a pure training time has reached kTPU*day(Tensor Processing Unit)[28]. Such a traversal method is not acceptable as the datasets expand and the participated modes increase. Therefore, an effective alternative method is needed to reduce the time cost and promote feature fusion.

In practice, We model the problem as an optimization problem. (As shown in Eq.2) However, due to the limitation of the gradient descent, it is challenging to ensure getting the extreme point[6, 16]. Since this training method takes a long time and it is difficult to ensure the results obtained, we turn to look for the upper bound of this extreme value and verify through experiments that our method is positively correlated to the traditional method. Meanwhile, our method greatly reduces the selecting time.(Meanwhile, this method proposed by us reduces the time needed for selection.)

In the step of multimodal feature fusion, the existing multimodal WS-TAL task[19, 33–35, 37] mainly uses various Transformer variants, such as Hybrid Attention Network(HAN), to fuse features, which implicitly assumes that feature similarity is equivalent to feature semantic similarity. This assumption is valid for single modes, but the differences of the similarities between modes are much greater in multi-mode problems.

To verify above judgment, the similarities within and between the extracted audio features and video features are calculated on the LLP dataset(Look, Listen, and Parse)[33].

Within one mode, 97 % of the two features with the maximum similarity were labeled identically,while between modes, this rate drop to less than 30 %. (See Table 1 for details) The probability reveals that there is a severe difference between feature similarity and semantic similarity in multimodal tasks.(Figure 1 shows one example) Therefore, It is inappropriate to utilize the multimodal Transformer structure directly for single-mode tasks, which will greatly damage the performance of feature fusion. Moreover, semantic constraints and corrections should be completed before inter-modal features before using them.

Also, it should be noted that, the difference in image features between adjacent frames is slight. However,as for TAL domain problems, most of the valuable information is hidden in these differences. The AVVP task will be adequately addressed if we can extract the information.

In light of these, referring to the idea of contrastive learning: reduce the distance between positive samples and increase the distance between negative samples, we introduce a multimodal contrastive loss. Due to the lack of segment level labels, we consider two modal features at the same time as positive samples, and those at different time as negative samples. It immensely enhances the correlation between feature similarity and feature semantic similarity among various modes, which is of great benefit for feature fusion. In this respect, Zhang[43] is similar to us. The differences between him and us lie in: 1. In his work, contrastive learning is introduced to solve the unimodal WS-TAL problem. He believes that(His reason is that) indistinguishable snippets could be easily misclassified and hurt the localization. 2. Different from him, we introduce our method in the multimodal task to solve the disalignment between the feature similarity and feature semantic similarity.

The main contributions of this paper are:

1. Pioneeringly,in the step of feature embedding, considering that feature extraction plays a vital role in feature fusion in multimodal tasks, we propose an effective method to select the collocation of different model feature extraction networks. Compared with the traversal method currently used, our method ensures effectiveness and reduces time complexity.

2. In the step of feature fusion, a multi-mode time-series contrastive loss(MTSC) is proposed to promote feature fusion between modes given the disalignment of feature similarity and semantic similarity in AVVP. It also helps to the expression of the difference information.

3. Both methods we proposed helps to facilitate the feature fusion. We demonstrated the feasibility of our method in the specific AVVP task. However,the two problems we proposed are common across the multimodal tasks. The solution can be generalized to other multimodal task.

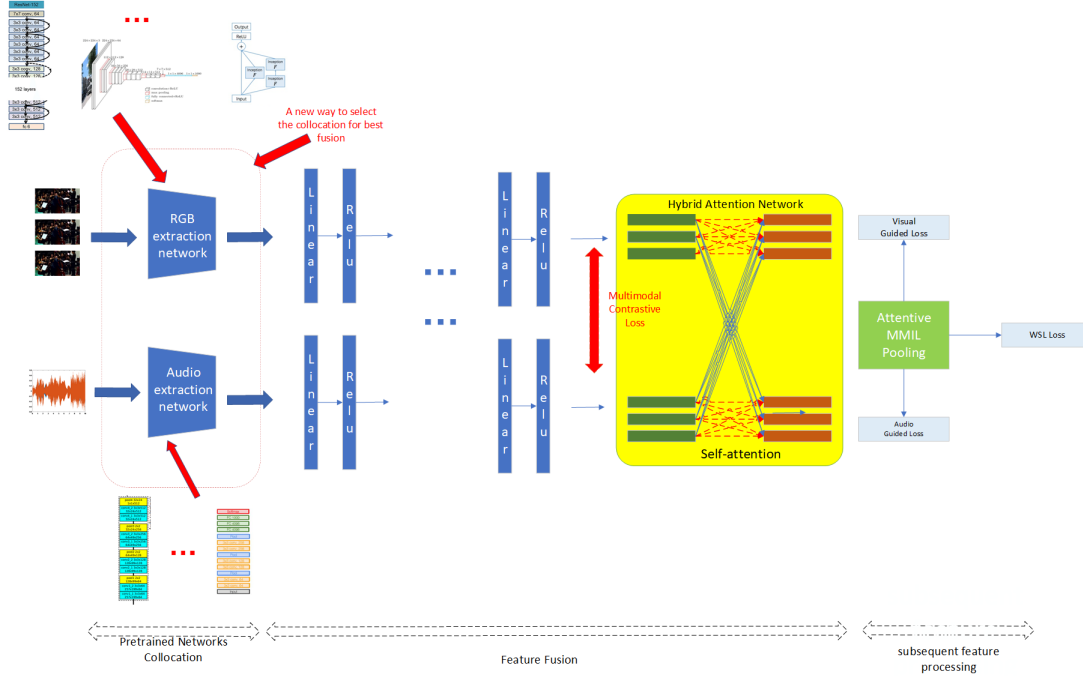


Figure 2. Our methods on baseline

2. Related works

1. **WS-TAL**: Weakly supervised temporal action localization only requires video-level annotation, which has attracted wide attention. *UntrimmedNets*[36] solve this problem by first classifying clip suggestions and then selecting relevant clips in a soft or hard manner. *STPN*[21] imposes a sparsity constraint to enforce the sparsity of the selected segment. *Hide-and-see* [27] and *MAAN*[41] extend the identification area by randomly hiding plaques or inhibiting the dominant response, respectively. *Zhong et al.*[47] introduced a progressive generation program to achieve a similar purpose. *W-TALC*[23] applies depth measurement learning to complement the multi-instance learning formula. Unlike actions in TAL, video events in audio-visual video parsing might contain motionless or even out-of-screen sound sources and the events can be perceived by either audio or visual modalities. .

2. **Contrastive Learning**: Contrastive Learning is a kind of self-supervised learning, which aims to learn knowledge by oneself from unlabeled data instead of relying on labeled data. Contrastive means learning to use internal data patterns to learn an embedding space in which correlated signals are aggregated, and non-correlated signals are distinguished by Noise Contrast Estimation (NCE)[12]. *CMC*(Contrastive Multiview Cocoding)[32] proposes a Contrastive learning framework that maximizes the mutual information between differ-

ent views of the same scenario to achieve view-invariant representation. *SIMCLR*[7] selects the negative sample by using an enhanced view of the other items in the small batch. Compared with the positive samples of comparison learning under the single mode, which also need to be obtained through data enhancement (*SIMCLR*) or different views of the same scene (*CMC*), timing multimode problems naturally have a pair of positive samples of different modes, and the information between different modes is supplementary rather than complementary. This naturally creates a unique application advantage for the application of comparative learning in multimode. To our knowledge, this paper is the first to combine contrastive learning with time-series information in multimodal TAL related tasks. The experimental results show that the contrast loss introduced in this paper is beneficial to multimodal time series localization.

3. **Multimodal feature fusion**: Fusion is a key research topic in multimodal studies, which integrates information extracted from different unimodal data sources into a single compact multimodal representation. There are three types of methods that are mainly used to fuse audio with image features, namely, simple operation-based[1, 22, 48], attention-based[2, 4, 8, 17, 25, 29, 38, 40], and tensor-based methods[10, 31, 42]. Among them, attention-based approach is the main research direction. However, attention needs the high correlation between feature similarity and feature semantic similarity, which is invalid between modes.

Therefore, a novel multimodal time-series contrastive loss is proposed by us to enhance the correlation.

4.Pre-trained feature extraction networks: For convenience, there are many pre-trained extraction networks[3, 9, 13, 14, 24, 26, 45, 46] for future embedding. The backbones and train sets of them vary, while all of the networks get a remarkable performance in multiple tasks after fine-tuning. However, in multimodal tasks, besides the remarkable performance of a single network, the mutual "understanding" between features in different modes extracted from the varied networks matters. We demonstrate experimentally that the latter is even more important than the former in multimodal tasks. Therefore, the collocation of pre-trained feature extraction networks with different modes is of great research value. Based on this, this paper proposes a novel method.

3. Approach

In this section, we introduce our method in detail. First, we explain some of the high frequency words(Section 3.1). Then, we describe our proposed method. In order, it consists two main parts: (i) a novel way to select the feature extraction networks(section 3.2). (ii) a multimodal time-series contrastive loss for better alignment.(section 3.3)

3.1. Explanation of key words

First of all, we explain some key words we often use. Like [44], we use the word "networks" to refer to pre-trained feature extraction networks; When we talk about AVVP's task network, we use the word "model". We use the words "mode" and "multimodal" respectively to describe the concepts of single and cross-modal.

3.2. A novel method for the collocation of pretrained feature extraction networks

3.2.1 Task objective function

Assume that there are N modes, and each mode has M pre-trained extraction network selections. Then the task's optimization goal is

$$\arg \max_{\Phi_1, \dots, \Phi_N} E(\Phi_1^{j_1}(x_1), \dots, \Phi_i^{j_i}(x_1), \dots, \Phi_N^{j_N}(x_1), \theta) \quad (1)$$

Here, Φ_i represents ith mode and Φ^{j_i} is any choice of one of the M networks in the ith mode, E is the target function, θ means parameters of the model.

Assume that each training time is T_0 and the inference time is T_1 , then the time to traversal all cases is $MN * (T_0 + T_1)$. Of course, the demand will consume more time in the actual situation because of time costs to adjust hyperparameters and other reasons, but it should be of this order.

The existing selection method is mainly traversal, and it is mathematically formulated as followed.

$$\max_E E(\phi(A_2), \theta), \arg \max_{\phi, \theta} E(\phi_i(A_1), \theta) \quad (2)$$

where $\phi_i \in \{\Phi_1^{j_1}(x_1), \Phi_2^{j_2}(x_1), \dots, \Phi_1^{j_i}(x_1), \dots, \Phi_N^{j_N}(x_1)\}$

However, this method has the disadvantages of a long time consuming and many traversal numbers.

3.2.2 A feasible solution

Since the training time cost is much longer than the test time cost, the most effective way to reduce the cost is to reduce the training time. Furthermore, due to the limitation of gradient descent, it is impossible to traverse every point on the loss hyperplane. So it is difficult to verify whether a comparison obtained by the above method is the same as the actual case (or just the result of inadequate tuning). In view of the above two purposes, we change the comparison object. To overcome the handicap of calculating the above Max formula's exact value, we turn to find an upper bound of the value, which is much easier to realize. Thus, we manage to change the comparative upper bound to measure the selection schemes quickly and effectively.

The traditional method is to find the optimal solution by training in the training set and testing in the test set. It directly calculates the optimal solution of all the specific tasks corresponding to the network, and then selects the optimal solution. This method is actually a joint solution of selecting the optimal model collocation and finding its. If we can skip the problem of finding the optimal model collocation result, choose the optimal solution in another way, and then calculate the performance of the optimal solution, then the costs can be reduced to $O(T_0+T_1)$. For multi-modal problems, this substitution method exists:

$$\left[\max_E E(\phi(A_2), \theta) \right]_{\arg \min_{\phi, \theta} L(\phi(A_1), \theta)} \quad (3)$$

$$= \max_E E(\phi(A_2), \theta) \arg \max_E E(\phi(A_2), \theta)$$

where $\lceil f \rceil$ represents the upper bound on the function f ; A_2 presents test set; A_1 presents training set and eval set; $\max_E E(\phi(A_2), \theta) \arg \max_E E(\phi(A_1), \theta)$ represents the process of training on the training set and optimizing on the eval set

A reasonable explanation is that: a model trained on the training set normally performs worse than one trained on the test set when both are tested on the test set – because the latter will have serious overfitting. Nevertheless, the overfitting is what we wanted. The more overfitting a collocation has, the easier it catch the multimodal informations. Therefore,

the latter is sufficient to compare the effects of different feature network options on feature fusion. In this way, we successfully separated the selection of the appropriate feature network collocation from the training of the entire network on the training set, reducing the time to $O(T1)$.

We regard the test result trained on the test set as F1, and the result trained on the training set as F2. The only drawback of the above argument is that it's hard to ensure the trend of change in F1 scores can reveal the trend of F2 when the collocation changes. In the subsequent experiments, we have proved that the change rule of F1 and F2 is roughly the same under various collocations of the mainstream feature extraction network. Such a result is understandable because the difference is mainly related to the subsequent model, whereas this model is fixed in our experiment.

3.3. MTSC

In view of the mismatch between feature similarity and feature semantic similarity in Multimodal tasks, as well as the fact that useful information is not effectively extracted in some modes, we refer to the idea of SIMCLR[7] and introduce Multimodal Time-series Contrastive Loss(MTSC).

In SIMCLR[7](As shown in Figure.3), for N samples in a minibatch, SIMCLR adopts two different data enhancement methods to obtain $2N$ sample points. SimCLR does not clearly indicate the negative sample pair. On the contrary, for a pair of positive sample points enhanced from a single data, The remaining $2(N-1)$ enhancement samples in minibatch are regarded as negative samples.

$$l_i = -\log \frac{\exp(\text{sim}(z_i^1, z_i^2)/\tau)}{\sum_{k=1}^N 1_{[k \neq i]} \exp(\text{sim}(z_i^1, z_k^2)/\tau)} \quad (4)$$

Different from the simple discrimination of positive and negative samples in SIMCLR, we have a priori information of time series in this task, that is, based on such a fact : semantic correlation decreases with the increase of time difference, so our loss function can be optimized as

$$l_i = -\log \frac{\exp(\text{sim}(\phi_i^1, \phi_i^2)/\tau)}{\sum_{k=1}^N 1_{[k \neq i]} \exp(T^* \text{sim}(\phi_i^1, \phi_k^2)/\tau)}, \quad (5)$$

$$T = F(k - i)$$

Due to the insensitivity of Eq.5 to the exponential part, the equation is modified as:

$$l = \frac{1}{N^2} \sum_{j=1}^N \sum_{i=1}^N [\text{sim}(\phi_i^1, \phi_j^2) - F(j - i) \cdot S_{i,j}] \quad (6)$$

$$S_{i,j} = T_{i,j} \cdot \text{sim}_{detach}(\phi_i^1, \phi_j^1)$$

where $T_{i,j}$ is a time lag coefficient which increases as the time lag decreases; sim_{detach} represents a detach of sim . We set $T_{i,j} = 1/|i - j|^5$.

Compared with the SIMCLR task, in the AVVP task, it is difficult to determine whether the features at different moments comprise a positive sample pair or a negative sample pair(As shown in Figure.4). Therefore, we choose to use the relationship between features in a single mode to make adaptive discrimination, namely sim_{detach} in the Eq.6.

We use time series information to increase the similarity of features of different modes, and reduce the similarity of features of different modes. While narrowing the gap between feature similarity and feature semantic similarity, the information contained in the feature difference value is re-extracted.

Algorithm 1 MTSC's main learning algorithm

Require: batch size N , structure of f, t, g, φ f is the feature extraction networks; t is part of the model before HAN; g is a Linear layer.

```

1: for sampled minibatch  $\{\mathbf{x}_k\}_{k=1}^N$  do
2:   for all  $k \in \{1, \dots, N\}$  do draw two modes  $\Phi \sim \varphi, \Phi' \sim \varphi$ 
3:     //the first mode
4:      $\tilde{x}_k = f_1(\Phi(x_k))$ 
5:      $h_k = t(\tilde{x}_k)$ 
6:      $z_k = g(h_k)$ 
7:     //the second mode
8:      $\tilde{x}'_k = f_2(\Phi'(x_k))$ 
9:      $h'_k = t(\tilde{x}'_k)$ 
10:     $z'_k = g(h'_k)$ 
11:   end for
12:   for all  $i \in \{1, \dots, N\}$  and  $j \in \{1, \dots, N\}$  do
13:      $s_{i,j} = z_i z'_j / (\|z_i\| \|z'_j\|)$ 
14:   end for
15:   define  $l_{i,j}$  as Eq. 6
16:    $L = \frac{1}{N^2} \sum_{j=1}^N \sum_{i=1}^N l(i, j)$ 
17:   update  $t, g$  to minimize  $L$ 
18: end for
return model  $t(\cdot)$ , and throw away  $g(\cdot)$ 

```

Like SimCLR, Algorithm 1 summarizes the proposed method.

In the specific application, to retain more semantic integrity of feature, Algorithm 1 could be integrated into the model training. To alleviate the problem of under-fitting states of different modes in the original model, we preserve the $g(\cdot)$ in visual brach.

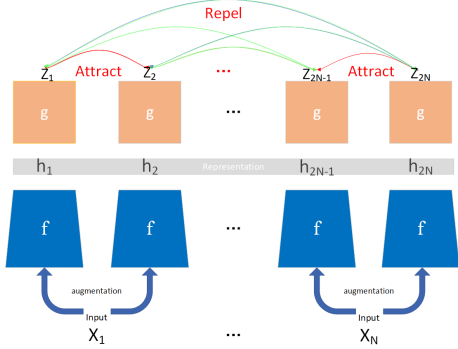


Figure 3. SimCLR[7]

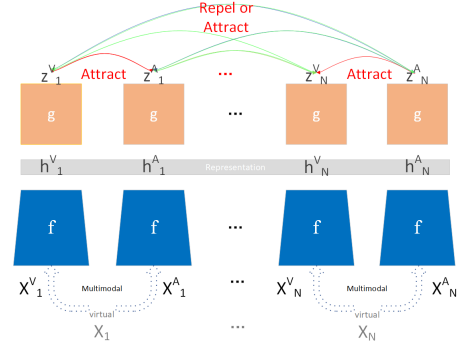


Figure 4. Multimodal SimCLR

4. Experiments

4.1. Implementation Details

We consider the same experimental settings as Tian did. For a 10-second-long video, we first sample video frames at 8fps, and each video is divided into non-overlapping snippets of the same length with eight frames in one second. Batch size and epochs are 16 and 40. The initial learning rate is $3E-4$ and will drop by multiplying 0.1 after every ten epochs. Our models optimized by Adam can be trained using one NVIDIA 2080 GPU.

Dataset we evaluate our method on the LLP datasets which is designed for the AVVP task. LLP contains 11,849 YouTube video clips spanning over 25 categories for a total of 32.9 hours collected from AudioSet [11]. A wide range of video events (e.g., human speaking, singing, baby crying, dog barking, violin playing, and car running, and vacuum cleaning etc.) from diverse domains (e.g., human activities, animal activities, music performances, vehicle sounds, and domestic environments) are included in the dataset.

Pretrained feature extraction networks We use two famous audio feature extraction networks including vggish[11] and openl3-audio[9], and five popular RGB image feature extraction networks including resnet-152[13], vgg19bn[26], polynet[46], SENET-154[14] and openl3-image[9].

Baselines We compare our method with the existing AVVP method[33] proposed by Tian. Since our contribution does not involve network design, our experiment is carried out on the network of Baseline. (As shown in Figure 2)

Evaluation Metrics As Tian[33] did, we evaluate them on parsing all types of events (individual audio, visual, and audio-visual events) under both segment-level and event-level metrics. To evaluate overall audio-visual scene parsing performance, we also compute aggregated results, where Type@AV computes averaged audio, visual, and audio-visual event evaluation results and Event@AV computes the F-score considering all audio and visual events for each sample rather than directly averaging results from different

event types as the Type@AV. We use both segment-level and event-level F-scores [20] as metrics. The segment-level metric can evaluate snippet-wise event labeling performance. We extract events with concatenating consecutive positive snippets in the same event categories and compute the event-level F-score based on mIoU = 0.5 as the threshold for computing event-level F-score results. Besides, Considering that this is a Multimodal Multiple Instance Learning (MMIL) problem with both segment-level task and event-level task, we consider a new Average Scores averaging the segments' and events' scores mentioned above as one of the final criteria for evaluating the performance of the model.

4.2. Experimental Comparison

4.2.1 the disalignment between feature semantic similarity and feature similarity in different modes

First, we verify a severe misalignment of the existing model between different modes, that is, the misalignment between feature semantic similarity and feature similarity. The details are shown in Table 1. Specifically, within one mode, 97.14 % of the two features with the maximum similarity were labeled identically, while between modes, this rate dropped to 29.18 %.

Then, the matching degree of feature semantic similarity and feature similarity between modes is compared in the feature fusion stage with/without the contrastive loss proposed by us. The details are shown in Table 2. The effectiveness of our method is evaluated from four aspects: recall-top1, distinguish, precision and performance. Recall-top1 is designed for the possibility that the most similar visual-audio segment-level pair have the same labels. It represents the ability of the model to separate positive and negative sample pairs; distinguish is used to measure the possibility that the most similar visual-audio segment-level pair have labels. It represents the ability of the model to distinguish between a background segment (without labels) and a target segment; precision represents the proportion of

Table 1. the possibility that the most similar segment-level pair that have the same labels.

	unimodal	multimodal
recall-top1	0.9714	0.2918

Table 2. Evaluation of our multi contrastive loss using existing multi-label classification metrics and our proposed classification metrics on LLP dataset. Part contrastive loss only restrict the features at the same time.

	raw	part contrastive loss	entire contrastive loss
recall	0.2918	0.0669	0.4112
distinguish	0.7193	0.9012	0.7973
precision	0.3186	0.2554	0.3538
Type@AV	0.533	0.543	0.554

M most similar video segments in each audio segment that have matching labels. M varies from segment to segment. It is the number of all video segments in the test set that match the label of each audio segment. In this setting, precision and recall have the same rate. Precision is the most accurate evaluation of the model’s ability to fit data sets. The above metric measures the ability of the part model before the transformer variants, and Type@AV represents the performance of the whole model.

As illustrated in the experiment, the misalignment between semantic similarity and feature similarity is more severe in multi-mode than in single-mode. In this case, the direct use of the Transformer is not appropriate. By comparing the results with and without contrastive loss, it is proved that our multimodal contrastive time-series loss is beneficial to alleviate this problem.

As shown in Table 2, our loss is divided into two parts: the part that restricts the same-time features; the part that restricts the different-time features. Both parts of our loss promote the final performance. But it’s interesting that when we only use the part that restricts the same-time features, it decreases the precision while highly improve the ability to separate the target segments from the background segments. It reveals two directions for optimizing baseline: one is to improve the ability to distinguish between target and background segments. The other is to improve the ability to distinguish between different target segments. An improvement in either of these two aspects will optimize model performance. Our losses contribute to the improvement of the model in both aspects.

4.2.2 The collocation of pre-trained networks

Secondly, We compare the results of our proposed method with the traditional method in networks collocation.

As shown in Figure5, We compared VGGISH and

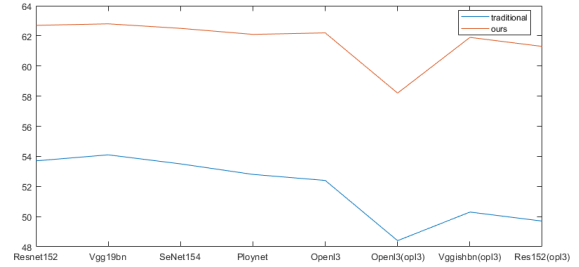


Figure 5. The comparison between the traditional traversal method and our method in evaluating the collocation of various feature extraction networks. The horizontal axis is the collocations of each mode’s pre-trained feature extraction networks, and the vertical axis is the evaluation index. The blue line represents the traditional method, and the orange line represents our method.

Table 3. The fine-tuning visual pretrained feature extraction networks test on the labeled images from LLP dataset.

	VGG19_bn	Resnet152	SENet-154	Polynet	openl3
accuracy	59.5	62.1	62.6	69.3	60.1

Table 4. The accuracy of distinguishing features of the same original data extracted by different feature extraction networks. We used pre-trained Resnet101 as the discriminator.

	Resnet	VGG19_bn	SeNet
PolyNet	99.9	99.9	99.8
SeNet	99.9	99.9	
VGG19_bn	99.8		

OpenL3 as the audio feature extraction networks combined with different image feature extraction networks as feasible solutions. (We did this because audio feature extraction has a very narrow range of network choices than image feature extraction, and this is where we think the Visual Audio task needs to be tackled.)

The experimental results verify that our method is consistent with the traditional traversal method. But our methods can reduce the time complexity.

In addition, an image dataset is created by collecting all the labeled frames from the LLP dataset and compare the classification accuracy of the fine-tuned image feature extraction networks. The results are shown in 3.

The pre-trained feature extraction network influences feature fusion(as shown in Figure 5) from two aspects: one is the ability of the network itself to extract the modal data information(Table 3), and the other is the ability of mutual fusion between the extracted features and other modal features. We call the latter the ability of mutual understanding among the pre-trained networks. By comparing Table 3 and Figure 5, we found that some networks with high classification accuracy were not as effective as those with low classification accuracy in Figure 5. The importance of feature

network collocation for feature fusion is proved.

Furthermore, we observe that the model gets the optimal result when the feature extraction models of different modes are similar. It is not hard to understand that besides the difference between source data characteristics in different modes, the network itself may also affect the result. As shown in Table 4, it is easy to distinguish the extraction results of the same input with different pre-trained networks, which indicates that the pre-trained extraction features have strong characteristics of the network itself. So when the pre-trained networks are similar, the extracted features are easier to be integrated in semantic. Therefore, the above overfitting phenomenon can be better understood as that when the structure of feature extraction networks of two modes is similar, the extracted features are easier to "understand" each other, so it is easier to carry out feature fusion.

4.2.3 Comparison with the SOTA

Finally, we compared our results with the baseline, and the results are shown in Table 5, proving that our results are superior to SOTA.

However, this is not the keypoint of our work. The central importance of our work lies in that we explore the specific manifestation of two problems common in multimodal tasks in the AVVP problem. The result better than the SOTA proves the existence and value of the methods we proposed.

5. Conclusion and future work

In this work, we propose two problems existing in the current research tasks of AVVP, 1. Multimodal pre-training model selection problem 2. Feature similarity and feature semantic similarity mismatch between modes. These problems exist not only in this task but also in many multimodal problems. We propose a fast and effective equivalent algorithm that greatly reduces the time spent on the first problem. For the second question, we are the first to propose a contrastive loss to constrain feature expression between modes in the TAL field, improving the feature fusion. The problems addressed by these two methods exist in the research field of AVVP but also are common in many multimodal task fields. Therefore, these two methods are of significant expansion and reference value.

6. Acknowledgements

Funding: This work was supported by the Ministry of Science and Technology of the People's Republic of China [grant number 2019YFC1511404]; and the National Natural Science Foundation of China [grant number 62002026].

References

- [1] A. Anastasopoulos, S. Kumar, and H. Liao. Neural language modeling with visual features. *arXiv preprint arXiv:1903.02930*, 2019.
- [2] P. Anderson, X. He, C. Buehler, D. Teney, M. Johnson, S. Gould, and L. Zhang. Bottom-up and top-down attention for image captioning and visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6077–6086, 2018.
- [3] R. Arandjelovic and A. Zisserman. Look, listen and learn. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 609–617, 2017.
- [4] D. Bahdanau, K. Cho, and Y. Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- [5] T. Baltrušaitis, C. Ahuja, and L.-P. Morency. Multi-modal machine learning: A survey and taxonomy. *IEEE transactions on pattern analysis and machine intelligence*, 41(2):423–443, 2018.
- [6] L. Bottou. Stochastic gradient descent tricks. In *Neural networks: Tricks of the trade*, pages 421–436. Springer, 2012.
- [7] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020.
- [8] J. Cho, J. Lu, D. Schwenk, H. Hajishirzi, and A. Kembhavi. X-ixmert: Paint, caption and answer questions with multi-modal transformers. *arXiv preprint arXiv:2009.11278*, 2020.
- [9] J. Cramer, H.-H. Wu, J. Salamon, and J. P. Bello. Look, listen, and learn more: Design choices for deep audio embeddings. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3852–3856. IEEE, 2019.
- [10] Y. Gao, O. Beijbom, N. Zhang, and T. Darrell. Compact bilinear pooling. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 317–326, 2016.
- [11] J. F. Gemmeke, D. P. Ellis, D. Freedman, A. Jansen, W. Lawrence, R. C. Moore, M. Plakal, and M. Ritter. Audio set: An ontology and human-labeled dataset for audio events. In *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 776–780. IEEE, 2017.
- [12] M. Gutmann and A. Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pages 297–304. JMLR Workshop and Conference Proceedings, 2010.
- [13] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [14] J. Hu, L. Shen, and G. Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7132–7141, 2018.
- [15] H. Huang, L. Su, D. Qi, N. Duan, E. Cui, T. Bharti, L. Zhang, L. Wang, J. Gao, B. Liu, et al. M3p: Learning universal representations via multitask multilingual multimodal pre-training. *arXiv preprint arXiv:2006.02635*, 2020.

Table 5. Comparisons with the state-of-the-art methods of the audio-visual video parsing task on the LLP test dataset. "contrastive loss partly" denotes that we only restrict the same-time features in our MSTC loss. "contrastive loss" denotes the whole proposed MSTC loss, "best network selection" denotes the method our proposed to select the best collocation of pre-trained feature extraction networks. "combination" denotes the combination of both methods we proposed.

Event type	Methods	Segment-level	Event-level
Audio	baseline	60.4	51.1
	+contrastive loss partly	60.6	52.0
	+contrastive loss	61.6	52.4
	best network selection	61.0	52.7
	combination	61.5	52.4
Visual	baseline	51.5	48.7
	+contrastive loss partly	52.7	48.5
	+contrastive loss	54.1	50.4
	best network selection	52.3	48.5
	combination	53.5	49.3
Audio-Visual	baseline	50.2	42.7
	+contrastive loss partly	50.4	43.2
	+contrastive loss	50.5	44.4
	best network selection	49.0	42.9
	combination	49.8	43.9
Type@AV	baseline	53.3	47.5
	+contrastive loss partly	54.1	47.9
	+contrastive loss	55.4	49.1
	best network selection	54.1	48.0
	combination	55.0	48.4
Event@AV	baseline	55.0	47.8
	+contrastive loss partly	55.1	48.0
	+contrastive loss	56.4	48.8
	best network selection	55.6	49.0
	combination	56.1	48.3

[16] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

[17] L. H. Li, M. Yatskar, D. Yin, C.-J. Hsieh, and K.-W. Chang. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019.

[18] Y. Li, Y. Pan, T. Yao, J. Chen, and T. Mei. Scheduled sampling in vision-language pretraining with decoupled encoder-decoder network. *arXiv preprint arXiv:2101.11562*, 2021.

[19] Y.-B. Lin, Y.-J. Li, and Y.-C. F. Wang. Dual-modality seq2seq network for audio-visual event localization. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2002–2006. IEEE, 2019.

[20] A. Mesaros, T. Heittola, and T. Virtanen. Metrics for polyphonic sound event detection. *Applied Sciences*, 6(6):162, 2016.

[21] P. Nguyen, T. Liu, G. Prasad, and B. Han. Weakly supervised action localization by sparse temporal pooling network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6752–6761, 2018.

[22] B. Nojavanasghari, D. Gopinath, J. Koushik, T. Baltrušaitis, and L.-P. Morency. Deep multimodal fusion for persuasive-ness prediction. In *Proceedings of the 18th ACM International Conference on Multimodal Interaction*, pages 284–288, 2016.

[23] S. Paul, S. Roy, and A. K. Roy-Chowdhury. W-talc: Weakly-supervised temporal activity localization and classification. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 563–579, 2018.

[24] H. Phan, L. Hertel, M. Maass, P. Koch, R. Mazur, and A. Mertins. Improved audio scene classification based on label-tree embeddings and convolutional neural networks. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 25(6):1278–1290, 2017.

[25] K. J. Shih, S. Singh, and D. Hoiem. Where to look: Focus regions for visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4613–4621, 2016.

[26] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.

[27] K. K. Singh and Y. J. Lee. Hide-and-seek: Forcing a network to be meticulous for weakly-supervised object and action localization. In *2017 IEEE international conference on*

- computer vision (ICCV)*, pages 3544–3553. IEEE, 2017.
- [28] W. Su, X. Zhu, Y. Cao, B. Li, L. Lu, F. Wei, and J. Dai. Vi-bert: Pre-training of generic visual-linguistic representations. *arXiv preprint arXiv:1908.08530*, 2019.
- [29] C. Sun, F. Baradel, K. Murphy, and C. Schmid. Learning video representations using contrastive bidirectional transformer. *arXiv preprint arXiv:1906.05743*, 2019.
- [30] C. Sun, A. Myers, C. Vondrick, K. Murphy, and C. Schmid. Videobert: A joint model for video and language representation learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7464–7473, 2019.
- [31] J. B. Tenenbaum and W. T. Freeman. Separating style and content with bilinear models. *Neural computation*, 12(6):1247–1283, 2000.
- [32] Y. Tian, D. Krishnan, and P. Isola. Contrastive multiview coding. *arXiv preprint arXiv:1906.05849*, 2019.
- [33] Y. Tian, D. Li, and C. Xu. Unified multisensory perception: weakly-supervised audio-visual video parsing. *arXiv preprint arXiv:2007.10558*, 2020.
- [34] Y. Tian, J. Shi, B. Li, Z. Duan, and C. Xu. Audio-visual event localization in unconstrained videos. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 247–263, 2018.
- [35] Y. Tian, J. Shi, B. Li, Z. Duan, and C. Xu. Audio-visual event localization in the wild. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition workshops*, 2019.
- [36] L. Wang, Y. Xiong, D. Lin, and L. Van Gool. Untrimmednets for weakly supervised action recognition and detection. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 4325–4334, 2017.
- [37] Y. Wu, L. Zhu, Y. Yan, and Y. Yang. Dual attention matching for audio-visual event localization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6292–6300, 2019.
- [38] H. Xu and K. Saenko. Ask, attend and answer: Exploring question-guided spatial attention for visual question answering. In *European Conference on Computer Vision*, pages 451–466. Springer, 2016.
- [39] Z. Yang, N. Garcia, C. Chu, M. Otani, Y. Nakashima, and H. Takemura. Bert representations for video question answering. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1556–1565, 2020.
- [40] Z. Yang, X. He, J. Gao, L. Deng, and A. Smola. Stacked attention networks for image question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 21–29, 2016.
- [41] Y. Yuan, Y. Lyu, X. Shen, I. W. Tsang, and D.-Y. Yeung. Marginalized average attentional network for weakly-supervised learning. *arXiv preprint arXiv:1905.08586*, 2019.
- [42] A. Zadeh, M. Chen, S. Poria, E. Cambria, and L.-P. Morency. Tensor fusion network for multimodal sentiment analysis. *arXiv preprint arXiv:1707.07250*, 2017.
- [43] C. Zhang, M. Cao, D. Yang, J. Chen, and Y. Zou. Cola: Weakly-supervised temporal action localization with snippet contrastive learning. *arXiv preprint arXiv:2103.16392*, 2021.
- [44] C. Zhang, Z. Yang, X. He, and L. Deng. Multimodal intelligence: Representation learning, information fusion, and applications. *IEEE Journal of Selected Topics in Signal Processing*, 14(3):478–493, 2020.
- [45] Q. Zhang, Z. Jiang, Q. Lu, J. Han, Z. Zeng, S.-H. Gao, and A. Men. Split to be slim: An overlooked redundancy in vanilla convolution. *arXiv preprint arXiv:2006.12085*, 2020.
- [46] X. Zhang, Z. Li, C. Change Loy, and D. Lin. Polynet: A pursuit of structural diversity in very deep networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 718–726, 2017.
- [47] J.-X. Zhong, N. Li, W. Kong, T. Zhang, T. H. Li, and G. Li. Step-by-step erasing, one-by-one collection: a weakly supervised temporal action detector. In *Proceedings of the 26th ACM international conference on Multimedia*, pages 35–44, 2018.
- [48] B. Zoph and Q. V. Le. Neural architecture search with reinforcement learning. *arXiv preprint arXiv:1611.01578*, 2016.