

Variational Autoencoder with CCA for Audio-Visual Cross-Modal Retrieval

JIWEI ZHANG*, Digital Content and Media Sciences Research Division, National Institute of Informatics, Japan

YI YU, Digital Content and Media Sciences Research Division, National Institute of Informatics, Japan

SUHUA TANG, Department of Computer and Network Engineering, Graduate School of Informatics and Engineering, The University of Electro-Communications, Japan

JIANMING WU, KDDI Research, Inc, Japan

WEI LI, School of Computer Science, Fudan University, China

Cross-modal retrieval is to utilize one modality as a query to retrieve data from another modality, which has become a popular topic in information retrieval, machine learning, and database. How to effectively measure the similarity between different modality data is the major challenge of cross-modal retrieval. Although several research works have calculated the correlation between different modality data via learning a common subspace representation, the encoder's ability to extract features from multi-modal information is not satisfactory. In this paper, we present a novel variational autoencoder (VAE) architecture for audio-visual cross-modal retrieval, by learning paired audio-visual correlation embedding and category correlation embedding as constraints to reinforce the mutuality of audio-visual information. On the one hand, audio encoder and visual encoder separately encode audio data and visual data into two different latent spaces. Further, two mutual latent spaces are respectively constructed by canonical correlation analysis (CCA). On the other hand, probabilistic modeling methods is used to deal with possible noise and missing information in the data. Additionally, in this way, the cross-modal discrepancy from intra-modal and inter-modal information are simultaneously eliminated in the joint embedding subspace. We conduct extensive experiments over two benchmark datasets. The experimental outcomes exhibit that the proposed architecture is effective in learning audio-visual correlation and is appreciably better than the existing cross-modal retrieval methods.

CCS Concepts: • **Information systems** → **Information extraction**.

Additional Key Words and Phrases: Cross-Modal Retrieval, Audio-Visual Correlation Learning

*Jiwei was involved in this work when he worked as an assistant researcher at the National Institute of Informatics, Tokyo, Japan.

Authors' addresses: Jiwei Zhang, Digital Content and Media Sciences Research Division, National Institute of Informatics, 2-1-2 Hitotsubashi, Chiyoda-ku, Tokyo, 101-8430, Japan, jiweizhang@nii.ac.jp; Yi Yu, Digital Content and Media Sciences Research Division, National Institute of Informatics, 2-1-2 Hitotsubashi, Chiyoda-ku, Tokyo, 101-8430, Japan, yiyu@nii.ac.jp; Suhua Tang, Department of Computer and Network Engineering, Graduate School of Informatics and Engineering, The University of Electro-Communications, 1-5-1 Chofugaoka, Chofu-shi, Tokyo, 182-8585, Japan, shtang@uec.ac.jp; Jianming Wu, KDDI Research, Inc, 2-1-15 Ohara, Fujimino, Saitama Prefecture, 356-8502, Japan, ji-wu@kddi-research.jp; Wei Li, School of Computer Science, Fudan University, 220 Handan Road, Yangpu District, Shanghai, Shanghai, 356-8502, China, weili-fudan@fudan.edu.cn.

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.

© 2018 Association for Computing Machinery.

0004-5411/2018/8-ART111 \$15.00

<https://doi.org/10.1145/1122445.1122456>

ACM Reference Format:

Jiwei Zhang, Yi Yu, Suhua Tang, Jianming Wu, and Wei Li. 2018. Variational Autoencoder with CCA for Audio-Visual Cross-Modal Retrieval. *J. ACM* 37, 4, Article 111 (August 2018), 20 pages. <https://doi.org/10.1145/1122445.1122456>

1 INTRODUCTION

Many services and applications connected with IoT involve more than one modality, which generate various data and information in different modalities (e.g., sensor, image, video, audio, text). Consequently, these multimodal data are accumulated over time at an unprecedented scale, which contain useful knowledge structures for describing various events and phenomena. This motivates us to use artificial intelligence to model the cognitive process of human with these big-scale data, and exploit it in cross-modal retrieval. Cross-modal representation learning is becoming more and more important for understanding the correlation among various modality data, which facilitates cross-modal retrieval between different modality data (such as text, audio, and visual). That is to say, using single modal information as a query is to retrieve the associated information of another modal. Searching results provided throughout several modalities can help customers to obtain comprehensive information about the desired event or topics.

Now, more and more researchers from academia and industry are starting to pay interest to cross-modal retrieval research. The major challenge of the cross-modal retrieval task is how to measure the content similarity between different modality data in the joint space, which is known as the heterogeneity gap. Many efforts have been carried out recently for cross-modal retrieval [1–5, 10–13, 17, 20, 45] between different modalities such as text-image [1], video-text [3], and audio-lyrics [5] with different levels of semantics.

Traditionally, linear projections are generally computed for measuring the correlation between different modalities. For example, Canonical Correlation Analysis (CCA) [15, 16, 18] is the most well-known learning methods. It finds the linear transformation of the two modalities of the input data by maximizing the pairwise correlation in the common subspace.

Benefiting from the speedy improvement of deep neural networks (DNN) technology, a lot of researchers have utilized DNN to capture nonlinear correlation and guide the model to learn the common representation subspace of multimodal data. Deep Canonical Correlation Analysis (DCCA) [17], is extensively utilized to learn complicated nonlinear transformations of different modalities. Some supervised deep cross-modal learning networks based on semantic category information are proposed to guide the model to learn extra discriminative representations. It aims to enable high-level semantic separation between different semantics in the common representation subspace. C-DCCA [19, 30] utilizes a deep learning network to extend the standard CCA, learning the nonlinear projection between pairs of different modality data while effectively preserving semantic information.

The success of present state-of-the-art methods has strongly relied on a giant quantity of high-quality labeled data. However, due to the wide range of information sources, it causes many noises and a lack of information. If machine learning methods ignore these uncertainties, it will inevitably lead to many meaningless statistical results. Therefore, it is necessary to use probabilistic modeling methods to deal with the uncertainty in the data. In MS-VAE [53], Zhu et al. presented a self-supervised architecture based on the VAE network (without semantic category information) to learn audio-visual correlation to realize cross-modal retrieval tasks.

Currently, existing methods have used semantic information to find out discriminative features from intra-modal or inter-modal data. However, in these cross-modal correlation learning methods, semantic category information has not been effectively utilized. Therefore, cross-modal representation learning still faces the following challenges:

- 1) How to effectively use the correlation between inter-modal and intra-modal data to learn the shared subspace representation better to enhance retrieval performance.
- 2) How to guarantee that the features encoded by the deep network can still reflect the structural distribution of the original data.
- 3) How to deal with the noise and missing data that may exist in the data, as the deep neural network ignores the uncertainty in the data.

To challenge these issues, in this paper, we propose a novel variational autoencoder (VAE) architecture for simultaneously mitigating the cross-modal discrepancy from intra-modal and inter-modal data. Probabilistic modeling of the distribution of input data is to generate latent variables that observe Gaussian unit distribution. In addition, it ensures that the feature extraction encoder can still well reflect the features of the original data. To obtain this goal, it minimizes the discriminative loss of samples in the semantic category subspace so as to supervise our architecture to learn discriminative features. In addition, the decoder is exploited to reconstruct latent variables. Therefore, each pairwise semantic information and the category information are entirely exploited to make sure that the learned representation is both discriminative in semantic structure and invariant across modalities.

Our proposed cross-modal retrieval architecture between audio and video has several significant contributions, as follows:

- i) The idea of VAE, while preserving the semantic structure of the original data, can also effectively learn the audio-visual correlation in the common representation subspace. Probabilistic modeling method is applied to avoid the inability to handle the possible noise and lacking information in the audio-visual data.
- ii) Audio-visual pairwise-level and category-level mutual latent spaces are separately constructed by CCA, which can enhance audio-visual correlation embedding and mitigate the cross-modal discrepancy between audio-visual data.
- iii) Several loss constraints are proposed to optimize the audio-visual representation learning in the mutual latent spaces, which can reinforce both audio-visual pairwise-level and category-level correlation.
- iv) Extensive experiments on benchmark datasets present that our approach is superior to the current state-of-the-art audio-visual cross-modal retrieval methods.

2 RELATED WORKS

Since different modalities usually have inconsistent distributions and feature representations, it is essential to learn a common space to bridge heterogeneous data and measure their correlation. The purpose is to project each modality's feature into a common representation subspace that can directly evaluate the correlation between different modalities. Various strategies have been proposed to analyze this common representation subspace. In this paper, we focus on real-valued representation learning to learn the common space to achieve audio-visual cross-modal retrieval. This category includes unsupervised approaches, semi-supervised approaches, and supervised approaches.

The unsupervised methods only use co-occurrence information to learn common representations for different modality data, such as CCA [18], has been widely utilized in cross-modal retrieval tasks via associating features of different modalities in a common representation space. Rasiwasia et al [14] proposed to maximize the correlation between visual features and text features in the common representation space via utilizing the CCA method. The methods of correspondence autoencoder (Corr-AE) [33] and deep canonically correlated autoencoder (DCCAE) [20] are representative works of this subcategory. In [53], the authors proposed to give visual or audio modality as input data to generate a pair of corresponding audio-visual data. The pair-based method is to learn

the correlation between samples via utilizing paired different modalities data. The authors in [41] suggest multi-view metric learning (MVML-GL) methods with global consistency and local smoothness.

Semi-supervised methods can exploit useful unlabeled data to improve supervised learning tasks when the labeled data are scarce or expensive. In Adaptively Unified Semi-supervised Learning (AUSL) for cross-modal retrieval in [6], the authors present a semi-supervised approach which uses both relevance of class labeled and unlabeled data to strengthen multi-modal correlation. In adaptive semi-supervised feature selection (ASFS) for cross-modal retrieval in [7], the authors learn potential information from unlabeled data to analyze the correlation between different modality data. In semi-supervised cross-modal retrieval (SSCMR) with label prediction in [8], the author applied the category prediction layer to guide the model to learn the common representation of different modalities, so as to achieve cross-modal retrieval tasks.

To learn more discriminative common representations, supervised methods exploit semantic category information to distinguish the samples from different categories. Supervised methods can be used to calculate the distance between samples of different categories, and the distance between samples of the same category should be as close as possible. In deep supervised cross-modal retrieval (DSCMR) [28], an approach of using semantic category information to learning discriminative features is proposed. In addition, to preserve cross-modal similarity, many existing cross-modal retrieval methods [32] guide the model to reduce the distance of each pair of modal data in the common Hamming space. In adversarial cross-modal retrieval (ACMR) [27], the category information is utilized to learn the discriminative features during the feature projection process.

This paper aims to fully apply both pairwise-level and category-level audio-visual information as constraints to guide the model to learn more discriminative and modal-invariant representations of different modalities data. Our proposed architecture bridges the heterogeneity gap while capturing discriminative features and ensuring that the representation extracted by the encoder still reflects the feature of the original data to enhance the accuracy of cross-modal retrieval. In particular, paired audio-visual and category-level information are exploited to simultaneously guide the model to learn more discriminative features. Probabilistic modelling methods are exploited to deal with the uncertainty in the data and make the decoder robust to noise. In addition, CCA is suggested to learn audio-visual correlation embedding in the mutual latent space. the idea of VAE is used to guarantee that the features extracted by the deep encoder can nevertheless reflect the structural distribution of the original data.

3 OUR METHOD

We first explain the motivation of cross-modal subspace representation based on VAE with discriminative loss function. Then we introduce the proposed architecture. Finally, we analyze the objective loss function and the detailed training process.

3.1 Motivations

There are various cross-modal retrieval approaches based on CCA, such as [49, 50]. However, these methods can only learn linear features. With the development of deep neural networks, Deep-CCA and its variants such as C-DCCA [19, 30] and TNN-C-CCA [29] are proposed to solve the non-linear correlation learning. However, in these existing methods, it is not considered whether the features extracted by the neural network encoder still retain the general structure of the original data. Therefore, in this work, the decoder is utilized in the architecture to reconstruct latent features from the deep convolutional encoder, which ensures that the underlying features can reflect the structure of the original data to improve retrieval performance. In addition, the deep neural network ignores the uncertainty in the data and cannot deal with the noise and missing information that

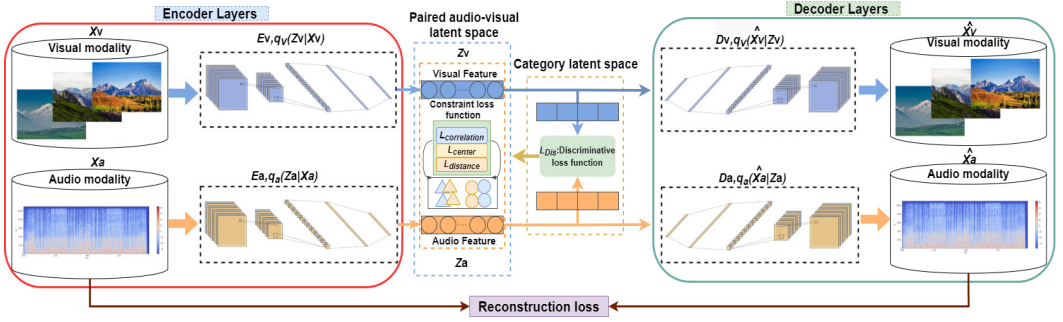


Fig. 1. The overall architecture of the proposed method. X_v and X_a are the input data from two modalities. \hat{X}_v and \hat{X}_a are the reconstruct data. E_v and E_a are deep convolutional encoders. Z_v and Z_a are latent representations from the outputs of E_v and E_a . D_v and D_a are decoders.

may exist in the data. To make the decoder robust to noise, Gaussian noise is added to the output of the encoder and probabilistic modeling methods are used to deal with the uncertainty in the data.

3.2 VAE-CCA Architecture

The proposed VAE with CCA architecture mainly consists of three parts as shown in Fig. 1: two branch encoder layers, mutual latent space, and two branch decoder layers.

3.2.1 Encoder layers. In our VAE-CCA model, we assume that two modality data are $X_v = \{x_v^i\}_{i=1}^m \in \mathbb{R}^{d_v \times m}$, $X_a = \{x_a^i\}_{i=1}^m \in \mathbb{R}^{d_a \times m}$, where m is the number of samples. d_v and d_a are the corresponding dimensions of the visual modality and the audio modality. We set deep convolutional encoders to three layers for two modalities. X_v and X_a are fed as inputs to deep convolutional encoders to process and we can obtain the latent representations $(Z_v | \theta_{e_v}) \in \mathbb{R}^{o \times m}$ and $(Z_a | \theta_{e_a}) \in \mathbb{R}^{o \times m}$, where θ_{e_v} , θ_{e_a} are respectively the parameters of deep convolutional encoders of visual and audio, and o is the output dimension of deep convolutional encoders.

3.2.2 Mutual latent space. Data from different modalities have different statistical characteristics and feature representations. Therefore, for cross-modal retrieval tasks, they cannot directly measure the correlation between different modalities data. In our work, we project the data of two modalities into common subspace, where we calculate the correlations between Z_v and Z_a with the following expression:

$$\arg \max_{\theta_{e_v}, \theta_{e_a}} \text{corr}(Z_a, Z_v) = \arg \max_{\theta_{e_v}, \theta_{e_a}} \frac{\text{cov}(Z_a, Z_v)}{\sqrt{\text{var}(Z_a)} \sqrt{\text{var}(Z_v)}}$$

where $\text{corr}(\cdot)$ is the correlation between Z_a and Z_v . $\text{cov}(Z_a, Z_v)$ is the covariance of Z_a and Z_v , and $\text{var}(Z_i)$ is the variance of Z_i , $i = a, v$.

3.2.3 Category latent space. The label matrix is denoted by $Y_i = [y_{i1}, y_{i2}, \dots, y_{ic}] \in \mathbb{R}^c$, where c is the total number of categories. If the i th sample belongs to the j th category, $y_{ij} = 1$, otherwise $y_{ij} = 0$, we convert the category into one-hot form.

3.2.4 Decoder layers. The decoder of each modality is composed of three layers of neural networks, which aims to reconstruct the latent representations from a shared wight layer and maintain the structural characteristics of the original data. We can obtain the outputs $(\hat{X}_v | \theta_{d_v})$ and $(\hat{X}_a | \theta_{d_a})$, where θ_{d_v} and θ_{d_a} are network parameters of the decoders.

3.3 Objective Loss Function Analysis

The goal of our architecture is to achieve the latent representation that can learn audio-visual correlation in the mutual space to mitigate the discrepancy between audio and visual data in different feature spaces. To do this, multiple loss functions are proposed to optimize our architecture.

3.3.1 VAE model Loss [38]. The objective function of our proposed architecture in reconstructing the data part is similar to the original VAE, where the log-probability $\log p(X)$ of the reconstruction data pair X_i is maximized from the desired mutual latent space Z_i , and i represents either the modality a (audio) or v (visual). We can get the variational lower bound of the objective function based on the data reconstruction of the original VAE network as follows:

$$\log p(X) \geq \mathbb{E}_{Z_i \sim q_i(Z_i | X_i)} [\log p(X | Z_i)] - \text{KL}(q_i(Z_i | X_i) \| p(Z_i)) \quad (1)$$

During the training process, the audio-visual pair from audio or visual input is reconstructed at each epoch. μ_i and σ_i for the Gaussian distribution are returned by the encoder, and Z_i is sampled from $\mathcal{N}(\mu_i, \sigma_i)$, $i = a, v$. In order to ensure the validity of the representation extracted by the encoder layer, we further add the decoders to separately reconstruct audio and visual data. The representations Z_v and Z_a from the shared weight layer are fed to the decoders and we can obtain the reconstruct data $(\widehat{X}_v | \theta_{d_v})$ and $(\widehat{X}_a | \theta_{d_a})$. Minimizing errors between $\widehat{X}_a, \widehat{X}_v$ reconstructed data and X_a, X_v original data can optimize the VAE networks. Therefore, the reconstruction loss for the network is:

$$\mathcal{L}_{rec} = \min_{\theta_{d_v}, \theta_{d_a}} \sum_{i=v,a} \left\| X_i - \widehat{X}_i \right\|_F^2 \quad (2)$$

We adopt the mean square error (MSE) [52] algorithm as the reconstruction loss function to calculate the audio-visual pair deviation between the reconstructed data \widehat{X}_i and the original data X_i . In the stage of training the VAE network to reconstruct the data, the total loss includes reconstruction loss and KL divergence.

$$\mathcal{L}_V = \mathcal{L}_{rec} + \mathcal{L}_{KL} \quad (3)$$

KL in the formula is used to represent the Kullback-Leibler divergence function, which is defined as: $\text{KL}(p(x) \| q(x)) = \int_x p(x) \log \frac{p(x)}{q(x)}$ is utilized to calculate the similarity between the distribution of audio-visual data.

3.3.2 \mathcal{L}_{corr} -correlation Loss [54]. In addition, we directly calculated the correlation of all paired audio-visual samples from the two modalities in the paired audio-visual latent space. As for the audio-visual pair data, we calculate the inter-modality and intra-modality correlation in the common representation subspace and then maximize the log-likelihood of the correlation:

$$\arg \max_{\theta_{e_v}, \theta_{e_a}} \text{inter-modality}_{discrimination} (Z_a, Z_v) =$$

$$\frac{1}{n^2} \sum_{a,v} (\log(1 + e^{\tau_{av}}) - \kappa_{av} \tau_{av})$$

where $\tau_{av} = \frac{1}{2} \text{corr}(Z_a, Z_v)$, $\kappa_{av} = 1 \{P_a, P_v\}$.

$$\arg \max_{\theta_{e_v}, \theta_{e_v}} \text{intra-modality}_{discrimination} (Z_v, Z_v) =$$

$$\frac{1}{n^2} \sum_{v,v} (\log(1 + e^{\gamma_{vv}}) - \zeta_{vv} \gamma_{vv})$$

where $\gamma_{vv} = \frac{1}{2} \text{corr}(Z_v, Z_v)$, $\zeta_{vv} = 1 \{P_v, P_v\}$.

$$\begin{aligned} \arg \max_{\theta_{e_a}, \theta_{e_a}} \text{intra-modality}_{discrimination}(Z_a, Z_a) = \\ \frac{1}{n^2} \sum_{a,a} \left(\log(1 + e^{\psi_{aa}}) - \phi_{aa} \psi_{aa} \right) \end{aligned}$$

where $\psi_{aa} = \frac{1}{2} \text{corr}(Z_a, Z_a)$, $\phi_{aa} = 1 \{P_a, P_a\}$. Generally speaking, our total cross-modal correlation loss function is defined as:

$$\begin{aligned} \mathcal{L}_{corr} = & \arg \max_{\theta_{e_v}, \theta_{e_a}} \text{inter-modality}_{discrimination}(Z_a, Z_v) \\ & + \arg \max_{\theta_{e_v}, \theta_{e_v}} \text{intra-modality}_{discrimination}(Z_v, Z_v) \\ & + \arg \max_{\theta_{e_a}, \theta_{e_a}} \text{intra-modality}_{discrimination}(Z_a, Z_a) \end{aligned} \quad (4)$$

$\text{corr}(\cdot)$ is utilized to calculate the similarity between audio-visual modality data, and $1\{\cdot\}$ is an indicator function. If the audio-visual modality data from the same category, its value is 1, otherwise it is 0.

3.3.3 \mathcal{L}_{dist} -distance Loss. We adopt minimizing the distance between all audio-visual pairs in the latent subspace representation to reduce the cross-modal heterogeneity difference between audio-visual modalities. Specifically, the Frobenius norm is used to directly measure the distance of all sample pairs in the paired audio-visual latent subspace. Finally, we define the modal invariance loss formula as follows:

$$\mathcal{L}_{dist} = \frac{1}{n} \|Z_v - Z_a\|_F \quad (5)$$

3.3.4 \mathcal{L}_{discr} -discriminative Loss [28]. To maintain the distinction of audio-visual samples from different categories after the features are mapped to the common space, we apply a simple linear layer classifier to predict the category of audio-visual modality samples projected in the latent space. Specifically, we connect a linear layer after the audio and visual encoder networks. The classifier makes use of the representation of the training data in the common subspace as input and generates a c -dimensional predictive category for each audio-visual sample. In the latent subspace of the semantic category, we adopt the Frobenius norm function to calculate the discriminative loss.

$$\mathcal{L}_{discr} = \frac{1}{n} \|P_a - y_{ai}\|_F + \frac{1}{n} \|P_v - y_{vi}\|_F \quad (6)$$

where $\|\cdot\|_F$ refers to the Frobenius norm, P_a and P_v are the projection matrix of the linear classifier, and n represents the batch size.

3.3.5 Center Loss [55]. In order to improve the discriminative power of deep learning features, we use an effective central loss function, which keeps the features of different categories separable while minimizing intra-category variations, as formulated in Equation 7.

$$\mathcal{L}_{center} = \frac{1}{2} \sum_{i=1}^m \|Z_v - c_{y_i}\|_2^2 + \frac{1}{2} \sum_{i=1}^m \|Z_a - c_{y_i}\|_2^2 \quad (7)$$

The $c_{y_i} \in \mathbb{R}^d$ refers to the y_i th category of the feature. The formulation effectively characterizes the intra-category variations.

According to analysis of the loss function for each part, the final objective function can be summarized as follows.

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{discr}} + \lambda_1 \mathcal{L}_V + \lambda_2 \mathcal{L}_{\text{corr}} + \lambda_3 \mathcal{L}_{\text{dist}} + \lambda_4 \mathcal{L}_{\text{center}} \quad (8)$$

In our experiments, according to the experimental results of parameters analysis, we set $\lambda_1 = 0.0001$, $\lambda_2 = 0.001$ and $\lambda_3 = 0.1$, and $\lambda_4 = 0.01$. We apply the gradient descent algorithm to optimize the objective loss function of the proposed architecture.

4 TRAINING STRATEGY

In the proposed architecture, we utilize two steps to train the architecture and optimize the whole network parameters.

4.1 Pre-training VAE Network

In the first step, we pre-train the VAE network utilizing Equation 3. We feed the cross-modal data X_a and X_v to deep convolutional encoders E_a and E_v , and acquire the reconstruction data \widehat{X}_a and \widehat{X}_v from the decoders D_a and D_v .

In the second step, We set the value of the learning rate to 3.5×10^{-4} . We adopt MSE to modify the objective function, where the error between the original data and the reconstructed data is minimized, the overall loss function is used to optimize the network and update encoders parameters θ_{e_a} , θ_{e_v} and decoders parameters θ_{d_a} , θ_{d_v} . The smaller the value of MSE is, the better the prediction performance of the proposed architecture is.

4.2 Training Entire Network

Finally, we train the entire network using Equation 8, minimizing the total loss including the discriminative loss, correlation loss, distance loss, center loss and the VAE network loss to update model parameters θ_{e_a} and θ_{e_v} , θ_{d_a} and θ_{d_v} . Algorithm 1 is the training process of our entire architecture.

5 EXPERIMENTS

In order to evaluate the performance of our proposed architecture, we conduct the experiment by comparing with ten remarkable baseline approaches on VEGAS and AVE datasets. Specifically, we first describe two datasets used in our paper, followed by the evaluation results.

ALGORITHM 1: The discriminative feature learning algorithm

Input: Visual samples for current batch: $X_v = \{x_v^1, \dots, x_v^n\}$;

Audio samples for current batch, $X_a = \{x_a^1, \dots, x_a^n\}$;

Corresponding labels for current batch, $Y = \{y_1, \dots, y_n\}$;

Output: The optimised parameters θ .

Initialize: λ_1 , λ_2 , λ_3 and λ_4 ; learning rate = 3.5×10^{-4} ; the number of iteration $t = 0$.

while *not converge* **do**

$t \leftarrow t + 1$

Train the VAE network using Equation 3

Optimize network parameters θ_{e_v} and θ_{e_a} of encoders and θ_{d_v} and θ_{d_a} of decoders

Renew the parameters of the sub-networks, by minising Equation 8.

end

5.1 Experiment Setup

5.1.1 Datasets Settings. Dataset with the audio-visual pairwise correlation and semantic category information are desired in the experiment. Therefore, we use the VEGAS [37] dataset and the Audio-Visual Event (AVE) [38] dataset.

The Audiovisual Events (AVE) data set is a subset of AudioSet [56], consisting of 28 event categories and a total of 4,143 videos. Each video has a period of 10s and utilizes audio-visual events as marker boundaries. The videos in the AVE dataset contain at least one audio-visual event with a duration of two seconds. The AVE dataset consists of various audiovisual event sounds in various fields such as human activities, animal activities, and music performances (for example, women talking, dog barking, playing guitar, etc.). In our experiment, we remove the music-related categories and keep 15 categories (Clock,motorcycle, train horn, bark, cat, bus, rodents/rats, toilet flush, acoustic guitar, frying, chainsaw, horse, helicopter, infant cry, truck)

The VEGAS dataset is a subset of the Google Audioset [56] with Amazon Mechanical Turk clean data, including 10 natural sound categories (chainsaw, helicopter, drum, printers, fireworks, dog, and etc.). The video duration is between 2-10 seconds, with an average of 7 second. It contains 28,103 videos and each video is marked by a single label. In our experiments, we used 5,621 videos for testing and 22,482 videos for training to evaluate our architecture. Finally, we summarize the statistical results of the two datasets in Table 1.

Table 1. Two benchmark data sets were used in our experiment, where N_{train} and N_{test} represent the number of training and testing audio-visual pairs, respectively. c represents the number of categories, D_{audio} and D_{visual} are the dimensions of the audio-visual features extracted by the VGGNet and VGGish networks, respectively.

Dataset	N_{train}	N_{test}	c	D_{visual}	D_{audio}
VEGAS	22482	5621	10	1024	128
AVE	1766	189	15	1024	128

5.1.2 Implementation details. In this work, our proposal architecture has two branch VAE networks, one of which is used for visual modality and the other is used for audio modality. The specific configuration of our proposed method is shown in Table 2. First, we utilize 19-layer VGGNet [51] to extract visual features of 512 dimensions, and a VGGish [39] to extract audio features of 128 dimensions. Then an encoder is applied in each sub-network to project different modality samples into a common subspace to learn the correlation of different modalities. A linear classifier layer is connected after each encoder to guide the network to obtain extra discriminative features. Finally, we reconstruct the learned latent representation into the original data dimension through the decoder. We apply GTX 2080 Ti GPU to train the proposed VAE network, and utilize the ADAM [36] to optimize the network parameter, the preliminary learning rate is set to 3.5×10^{-5} , and the maximum number of training times is set to five hundred.

5.1.3 Warmup Learning. The performance of the deep learning model is closely related to the setting of the learning rate. The proposal framework initially conducts model training with a large and constant learning rate. In the actual model training, first, we linearly increase the learning rate from 3.5×10^{-5} to 3.5×10^{-4} in the first 10 periods. Then, the learning rate decays to 3.5×10^{-5} in the 40th period, and is set to 3.5×10^{-6} in the 70th period and later.

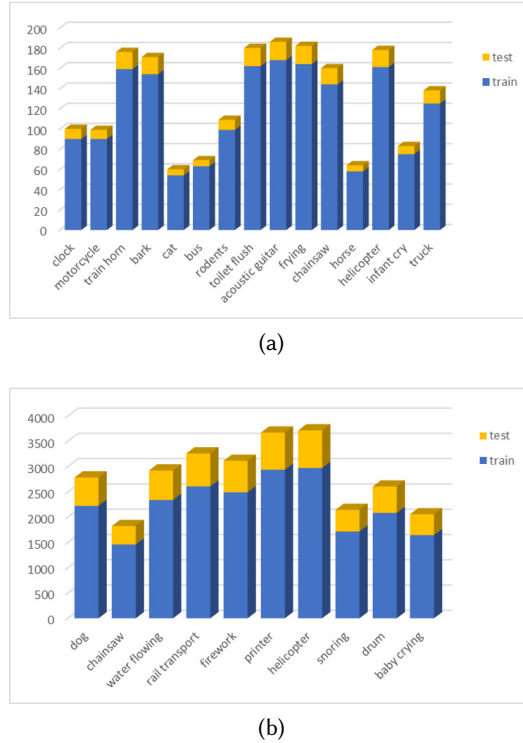


Fig. 2. The number of videos with different categories in training set and testing set in (a) AVE dataset, and (b) VEGAS dataset.

Table 2. Configuration of our proposed method.

Input	Visual Branch	Audio Branch
	1024-D	128-D
Layer 1 (Encoder)	512, fully-connected, linear	512, fully-connected, linear
Layer 2 (Common subspace layer)	64, fully-connected, linear	64, fully-connected, linear
Layer 3 (Semantic category layer)	10, fully-connected, linear	10, fully-connected, linear
Layer 4 (Decoder)	64, fully-connected, linear	64, fully-connected, linear
Layer 5 (Decoder)	512, fully-connected, linear	512, fully-connected, linear
Layer 6 (Decoder)	1024, fully-connected, linear	128, fully-connected, linear

5.2 Experimental Results

5.2.1 Evaluation Metric. We utilize the cosine value on the VEGAS and AVE datasets to calculate the retrieval similarity between different modalities, and adopt the average of all returned accuracy (mAP) as the evaluation metric. mAP is a performance evaluation standard widely used in cross-modal retrieval research [1, 21, 27]. It measures the ranking information and accuracy in joint consideration. In our experiment, we summary the mAP scores of two comparison methods for different cross-modal retrieval tasks:

- Retrieving audio samples using visual queries (Visual2Audio).

- Retrieving visual samples using audio queries (Audio2Visual).

5.2.2 *Comparison with Existing Approaches.* As shown in Table 3 and Fig. 3, we utilized the mAP metric and the PRC metric to record the experimental results of the audio-visual cross-modal retrieval on the VEGAS dataset, and visualized the experimental results. To confirm the superiority of our proposed architecture, we compare with eleven existing most advanced cross-modal retrieval methods, including three traditional methods: CCA [18], KCCA [35], and C-CCA [34], as well as eight deep learning-based methods, namely DCCA [17], C-DCCA [19, 30] UGACH [10], AGAH [25], UCAL [26], ACMR [27], DSCMR [28] and TNN-C-CCA [29].

Table 3. Comparison with Existing Approaches on VEGAS dataset in terms of mAP. The highest score is shown in boldface.

Method	Audio2Visual	Visual2Audio	Average
CCA	0.332	0.327	0.330
KCCA	0.288	0.273	0.281
DCCA	0.478	0.457	0.468
C-CCA	0.711	0.707	0.709
C-DCCA	0.722	0.716	0.719
UGACH	0.182	0.179	0.181
AGAH	0.578	0.568	0.573
UCAL	0.446	0.436	0.441
ACMR	0.465	0.442	0.454
DSCMR	0.732	0.721	0.727
TNN-C-CCA	0.751	0.738	0.745
Ours	0.811	0.813	0.812

Table 4. Comparison with Existing Approaches on AVE dataset in terms of mAP. The highest score is shown in boldface.

Method	Audio2Visual	Visual2Audio	Average
CCA	0.190	0.189	0.190
KCCA	0.133	0.135	0.134
DCCA	0.221	0.223	0.222
C-CCA	0.153	0.152	0.153
C-DCCA	0.230	0.227	0.229
UGACH	0.165	0.159	0.162
AGAH	0.200	0.196	0.198
UCAL	0.153	0.150	0.152
ACMR	0.162	0.159	0.161
DSCMR	0.314	0.256	0.285
TNN-C-CCA	0.253	0.258	0.256
Ours	0.358	0.343	0.350

• CCA [18] projects the features of different modalities into a common subspace, and realizes cross-modal retrieval tasks by maximizing the correlation between modal samples.

- *KCCA* [35] improves the CCA algorithm by introducing the concept of a "kernel trick" for learning common space representation. We utilize Gaussian kernel as the kernel function of CCA in the comparison experiment.

- *C – CCA* [34] (Cluster-CCA) clusters different modalities data into several categories, and try to enhance the correlation intra-cluster.

- *DCCA* [17] utilizes deep learning network to solve nonlinear projection problem, apply CCA-like objective function to maximize the correlation of different modal samples.

- *C – DCCA* in [19, 30] is a combination of Cluster-CCA and DCCA, which learn the non-linear representation by deep learning method into several related clusters to optimize the correlation.

- *UGACH* [10] utilizes GAN network to extract potential features of cross-modal data.

- *AGAH* [25] the multi-label attention module of adversarial learning applies to enhance the ability to distinguish between cross-modal representations.

- *UCAL* [26] maximizes the correlation between visual-text modalities in the common representation space. The classifier predicts the visual-text modalities of the learned features, and applies adversarial learning ideas to add regularization to the model.

- *ACMR* [27] adopts adversarial training ideas, utilizing classifiers and feature projections to guide the model to learn modal invariant and discriminative features.

- *DSCMR* [28] guides the network to learn the discriminative features in the category subspace and feature common subspace through the method of supervised learning.

- *TNN – C – CCA* [29] uses the triple loss function to reduce the category spacing of features in the common representation subspace, thereby improving the retrieval performance of Cluster-CCA.

The performances of our proposed method and comparison algorithms on AVE datasets are reported in Table 4. From the presented results, we can have the following observations:

- Our proposed architecture can achieve the best performance on the AVE datasets in terms of mAP, which verifies the impact of improved retrieval performance via the correlations between different modalities data.

- Our proposed architecture is extensively better than both CCA [18] and KCCA [35] among most cases, e.g., on the AVE dataset, CCA and KCCA are only 0.190 and 0.134 for mAP. That is because they are unsupervised learning methods, which do not consider discriminative label information of different modalities.

- The reason for the poor clustering overall performance of DCCA is that it can not reconstruct the information to make sure that the representation after encoding the network can nevertheless reflect the structure of the original data, which causes mAP to be 0.222 on the AVE dataset. Additionally, our proposed architecture also performs better than DSCMR and TNN-C-CCA because these two methods cannot make full use of the correlations among the inter-modal data.

5.2.3 Convergence Analysis. In this section, we train the proposed architecture on the VEGAS and AVE datasets, and use the output loss value to draw the change curve of the total loss value of the objective function to investigate the convergence of the proposed architecture. As shown in Fig. 4, the value of the total loss function decreases as the number of iterations increases, and the values approach to be a fixed value after a few iterations (less than 20 iterations), where each iteration includes 100 epochs. Therefore, our proposed optimization algorithm is reliable and converges quickly.

5.2.4 Parameters Analysis. In our architecture, there are four regularization parameters λ_1 , λ_2 , λ_3 and λ_4 , we simultaneously adjust them to obtain the best model, However, for simplicity and evaluating the effect of each parameter in our experiments, we fix one and vary the other three for each time.

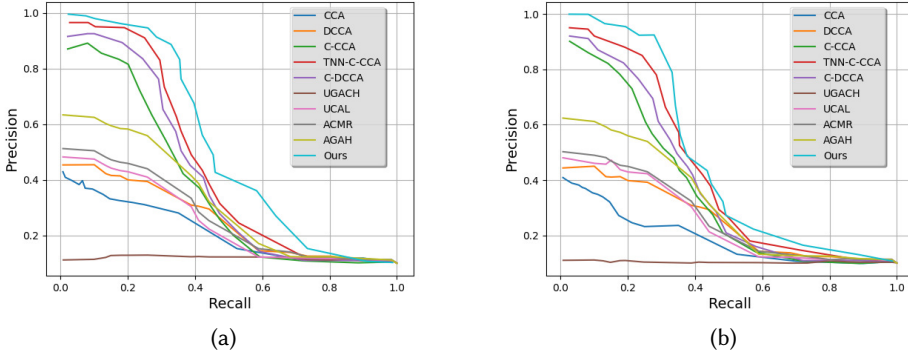


Fig. 3. We visualize the experimental values of PRC and compare the other nine different models on the VEGAS dataset. (a) is for audio2visual retrieval and (b) is for visual2audio retrieval.

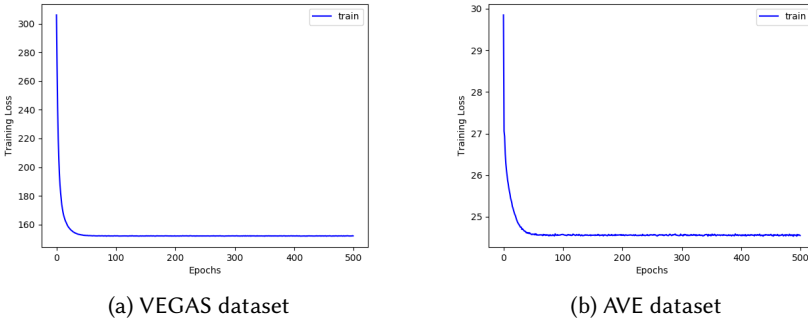


Fig. 4. We visualize the loss curve of our architecture on the VEGAS and AVE datasets. And set 500 epochs to train the entire network and obtain a loss value for each 100 epochs.

Firstly, we set the regularization parameters of the VAE calculation λ_1 and modify the regularization parameters of the correlation error, the distance error and the center error λ_2 , λ_3 and λ_4 in range $\{0.001, 0.01, 0.1, 1\}$. Then we fix λ_2 and also vary λ_3 and λ_4 in the same range. Since the approach of adjusting the parameters is the equal on both datasets, we solely exhibit the changes of the parameters in the VEGAS dataset. From the Fig. 5, we notice that:

- Our architecture can obtain the best mAP values on VEGAS dataset when $\lambda_1 = 0.0001$, $\lambda_2 = 0.001$, $\lambda_3 = 0.1$ and $\lambda_4 = 0.01$;
- In our architecture, when λ_3 and λ_4 are fixed, the change range of λ_1 and λ_2 is relatively large; when λ_1 and λ_2 are fixed, the change range of λ_3 and λ_4 is relatively small.

5.2.5 Role of Category. In addition, we investigate the effectiveness of common subspace representation via audio-visual retrieval tasks. Fig. 6 shows the average accuracy (AP) score of each category search after comparing our model with C-CCA, TNN-C-CCA and DSCMR models on the VEGAS dataset. It can be roughly seen in Fig. 6 that the higher the AP value is, the easier it is to retrieve, and the retrieval accuracy of different categories varies greatly. When the audio is utilized as a query, we

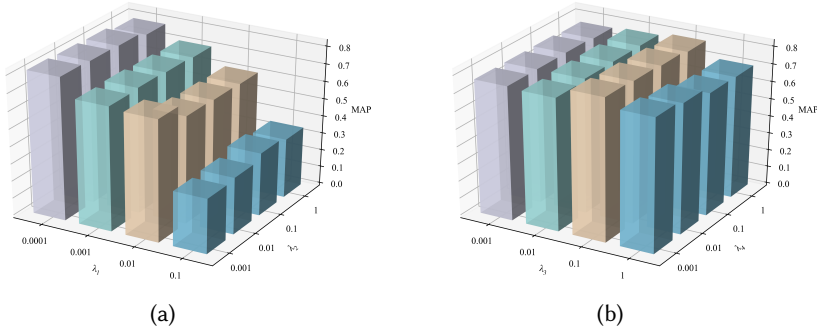


Fig. 5. The effect of parameters λ_1 , λ_2 , λ_3 and λ_4 on VEGAS dataset, where λ_1 is the regularization parameter of the VAE calculation, λ_2 is the regularization parameter of the correlation calculation, λ_3 is the regularization parameter of the latent loss error and λ_4 is the regularization parameter of the center calculation. (a) is the retrieval results in terms of mAP, when fixing λ_3 , λ_4 and varying λ_1 and λ_2 . (b) is retrieval results in terms of mAP, when fixing λ_1 , λ_2 , and varying λ_3 , and λ_4 .

can see that “baby crying” and “drum” use the proposed method to get the highest AP. When using the our proposed architecture (VAE-CCA) for retrieving the “helicopter”, audio2visual can reach 82.3%, while visual2audio is only 78.3%, a difference of nearly 4.0%. By comparing with methods such as DSCMR, TNN-C-CCA, and C-CCA, our proposed architecture can obtain retrieval accuracy higher than other methods in each category, which verifies the effectiveness of our architecture.

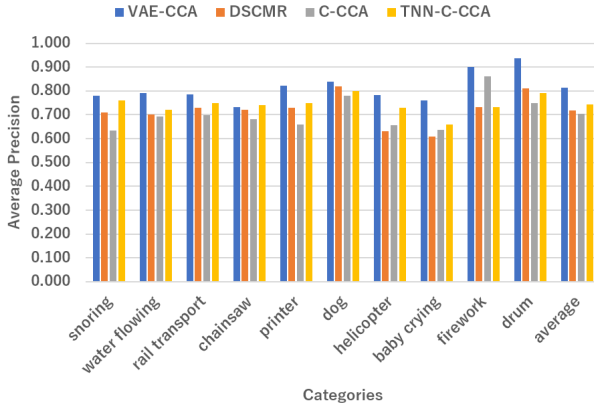
5.2.6 Ablation Study. The total loss function of our proposed architecture is composed of multiple loss functions to minimize the discriminative loss in the category latent space and the common representation subspace, and the distance loss and centre loss in the common representation subspace, respectively. We conducted ablation experiments on the VEGAS data set to examine the impact of the different loss functions on the whole performance of the proposed architecture. The experimental results are shown in the Table 5.

From Table 5 we can observe:

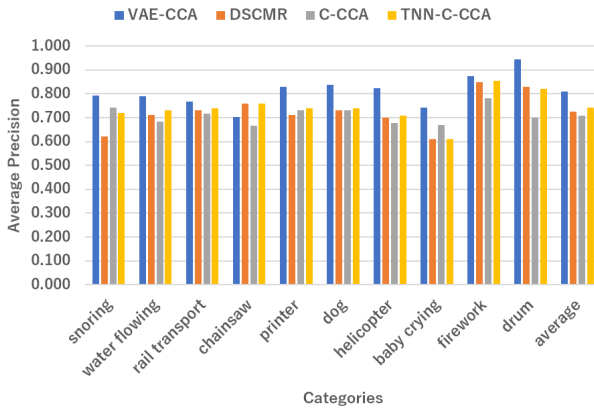
- Correlation constraints have a positive affect on retrieval performance. They maximize the correlation of the information between the modalities and obtain a better representation of the common subspace.
- The distance loss has an impact on on the proposed architecture. Based on the above experimental results, all the loss functions in our proposed architecture can be combined to obtain properly retrieval performance.

Table 5. Ablation study on VEGAS dataset in terms of mAP.

Methods	Audio2Visual	Visual2Audio	Average
With center loss	0.130	0.145	0.138
With correlation loss	0.520	0.404	0.462
With distance	0.651	0.638	0.644
Full our proposed method	0.811	0.813	0.812



(a)



(b)

Fig. 6. We compare the proposed method with the C-CCA, TNN-C-CCA, and DSCMR methods in the VEGAS dataset, and show the results of each category. According to the visualized statistical results, the retrieval accuracy rates of different categories are quite different. (a) visual2audio and in (b) audio2visual.

5.2.7 Visualization of the Learned Representation. Here we adopt the t-SNE [40] approach to map the common representation subspace of audio and visual samples into a two-dimensional visualization plane to learn about the effectiveness of the proposed architecture. Fig. 7(a)-(c) exhibit the distribution of original 1,024-dimensional visual features and 128-dimensional audio features. We can see that the distribution of visual modalities and audio modalities in the VEGAS Dataset is largely different, and the sample category spacing is small. It is not always easy to distinguish the sample categories effectively. Fig. 7(d)-(f) show the two-dimensional distributions of the visual and audio representations in the common subspace by C-CCA method. from Fig. 7(e)-(f) We can see that C-CCA embeds category information into the feature representation space, where clusters of different colors represent different categories. It can be observed that these clusters are not completely distinguishable. From the visualization results, the embedding of our VAE-CCA is

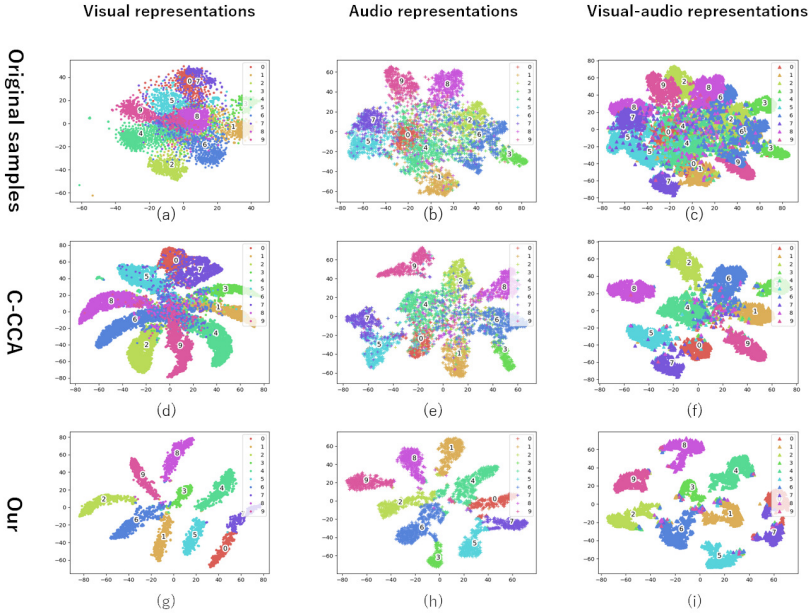


Fig. 7. The t-SNE method is used to visualize the audio-visual test data in the VEGAS dataset. The circles represent samples from visual and audio modality. We utilize the same color to represent samples belonging to the same semantic category, visualize the learned representation subspace via utilizing the t-SNE approach, showing the original data distribution of audio, visual, and audio-visual. audio-visual, and visualize the feature representation subspace learned by the C-CCA method and the feature representation subspace learned by our proposed method. The circle represents audio data, and the triangle sign represents visual-audio data.

much better than C-CCA embedding. The function of distinguishing loss in common subspace and semantic category subspace can distinguish samples from different semantic categories, and efficaciously divide the samples of joint space into independent semantic groups.

In addition, to verify the effectiveness of common subspace representation, we conducted audio2visual retrieval experiments. Comparing our model with the other three best models on the VEGAS dataset, Fig. 8 shows the audio2visual retrieve results of ACMR, AGAH, TNN-C-CCA, and our model. We apply the audio “chainsaw” as the query, and we can observe that the AP of our model is 83.2% of all ranking lists. For other models, the AP of the ACMR model is 41.6% of all rankings; the AP of the AGAH model is 56.3% of all rankings.

Finally, we adopt the confusion matrix evaluation metric on the VEGAS dataset to conduct video and audio cross-modal retrieval experiments to evaluate the accuracy of the proposed cross-modal retrieval architecture and then calculate the confusion matrix and visualized the results. As shown in Fig. 9, the number of retrieved samples is proportional to the color brightness. The horizontal axis represents the predicted category, and the vertical axis represents the actual category. The number of samples retrieved correctly is displayed on the diagonal, and the number retrieved incorrectly is displayed on the area outside the diagonal. From our experimental results, the correctly retrieved samples are concentrated on the diagonal, which verifies that our proposed architecture has better retrieval performance.

Audio query	Chainsaw					
	ACMR:AP=0.416	AGAH:AP=0.563	TNN-CCA:AP=0.801	Our:AP=0.832		
Top-5 retrieved visuals	chainsaw					
	water flowing					
	chainsaw					
	printer					
	firework					

Fig. 8. Visualize the results of audio2visual retrieval: compare our proposed architecture with the other three best existing methods TNN-CCA, AGAH and ACMR models. Apply audio as a query and display the top five retrieved visuals.

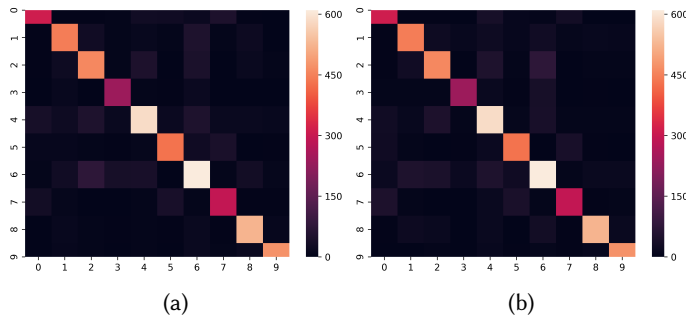


Fig. 9. The confusion matrix achieved on the VEGAS dataset with our proposed architecture. The figure (a) is for audio2visual retrieval, the figure (b) is for visual2audio retrieval.

6 CONCLUSION

In this paper, we propose a novel architecture based on VAE (VAE-CCA) to guide the model to learn more discriminative features via making full use of semantic category information and modal-invariant representations of different modalities data. We apply CCA to learn audio-visual correlation embedding in the mutual latent space. Our architecture can mitigate the discrepancy between audio-visual data while capturing discriminative features, and making sure that the representation processed by way of the encoder will keep the traits of the unique data. Since the wide range of datasets, there are problems such as noise and lack of information, which

lead to uncertainty in the data. In this paper, we utilize probabilistic modeling methods to deal with this issue. We carried out many comparative experiments on two benchmark datasets and comprehensively analyzed the experimental results to exhibit the effectiveness of the proposed cross-modal retrieval architecture.

In future research, we desire to apply our architecture to different kinds of cross-mode retrieval, such as visual2text, audio2text, and video2text. We will additionally consider to apply adversarial learning methods to enhance our cross-modal retrieval performance, and strive to extend our current architecture to achieve cross-modal data fusion to address the issue of insufficient data.

REFERENCES

- [1] K. Wang, Q. Yin, W. Wang, S. Wu, and L. Wang, "A comprehensive survey on cross-modal retrieval," *arXiv preprint arXiv:1607.06215*, 2016.
- [2] Y. Wei, Y. Zhao, C. Lu, S. Wei, L. Liu, Z. Zhu, and S. Yan, "Cross-modal retrieval with cnn visual features: A new baseline," *IEEE transactions on cybernetics*, vol. 47, no. 2, pp. 449–460, 2016.
- [3] X. Xu, L. He, H. Lu, L. Gao, and Y. Ji, "Deep adversarial metric learning for cross-modal retrieval," *World Wide Web*, vol. 22, no. 2, pp. 657–672, 2019.
- [4] Y. Yu, Z. Shen, and R. Zimmermann, "Automatic music soundtrack generation for outdoor videos from contextual sensor information," in *Proceedings of the 20th ACM international conference on Multimedia*, 2012, pp. 1377–1378.
- [5] Y. Yu, S. Tang, F. Raposo, and L. Chen, "Deep cross-modal correlation learning for audio and lyrics in music retrieval," *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, vol. 15, no. 1, pp. 1–16, 2019.
- [6] L. Zhang, B. Ma, J. He, G. Li, Q. Huang, and Q. Tian, "Adaptively unified semi-supervised learning for cross-modal retrieval," in *IJCAI*, 2017, pp. 3406–3412.
- [7] E. Yu, J. Sun, J. Li, X. Chang, X.-H. Han, and A. G. Hauptmann, "Adaptive semi-supervised feature selection for cross-modal retrieval," *IEEE Transactions on Multimedia*, vol. 21, no. 5, pp. 1276–1288, 2018.
- [8] D. Mandal, P. Rao, and S. Biswas, "Semi-supervised cross-modal retrieval with label prediction," *IEEE Transactions on Multimedia*, vol. 22, no. 9, pp. 2345–2353, 2019.
- [9] R. R. Shah, Y. Yu, and R. Zimmermann, "Advisor: Personalized video soundtrack recommendation by late fusion with heuristic rankings," in *Proceedings of the 22nd ACM international conference on Multimedia*, 2014, pp. 607–616.
- [10] J. Zhang, Y. Peng, and M. Yuan, "Unsupervised generative adversarial cross-modal hashing," in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [11] K. Wang, R. He, L. Wang, W. Wang, and T. Tan, "Joint feature selection and subspace learning for cross-modal retrieval," *IEEE transactions on pattern analysis and machine intelligence*, vol. 38, no. 10, pp. 2010–2023, 2015.
- [12] L. Zhang, B. Ma, G. Li, Q. Huang, and Q. Tian, "Generalized semi-supervised and structured subspace learning for cross-modal retrieval," *IEEE Transactions on Multimedia*, vol. 20, no. 1, pp. 128–141, 2017.
- [13] J. Wu, Z. Lin, and H. Zha, "Joint latent subspace learning and regression for cross-modal retrieval," in *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2017, pp. 917–920.
- [14] N. Rasiwasia, J. Costa Pereira, E. Coviello, G. Doyle, G. R. Lanckriet, R. Levy, and N. Vasconcelos, "A new approach to cross-modal multimedia retrieval," in *Proceedings of the 18th ACM international conference on Multimedia*, 2010, pp. 251–260.
- [15] D. R. Hardoon, S. Szedmak, and J. Shawe-Taylor, "Canonical correlation analysis: An overview with application to learning methods," *Neural computation*, vol. 16, no. 12, pp. 2639–2664, 2004.
- [16] S. E. Leurgans, R. A. Moyeed, and B. W. Silverman, "Canonical correlation analysis when the data are curves," *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 55, no. 3, pp. 725–740, 1993.
- [17] G. Andrew, R. Arora, J. Bilmes, and K. Livescu, "Deep canonical correlation analysis," in *International conference on machine learning*. PMLR, 2013, pp. 1247–1255.
- [18] H. Hotelling, "Relations between two sets of variates," in *Breakthroughs in statistics*. Springer, 1992, pp. 162–190.
- [19] Y. Yu, S. Tang, K. Aizawa, and A. Aizawa, "Category-based deep cca for fine-grained venue discovery from multimodal data," *IEEE transactions on neural networks and learning systems*, vol. 30, no. 4, pp. 1250–1258, 2018.
- [20] W. Wang, R. Arora, K. Livescu, and J. Bilmes, "On deep multi-view representation learning," in *International conference on machine learning*. PMLR, 2015, pp. 1083–1092.
- [21] Y. Peng, X. Huang, and J. Qi, "Cross-media shared representation by hierarchical learning with multiple deep networks," in *IJCAI*, 2016, pp. 3846–3853.
- [22] Y. Peng, J. Qi, X. Huang, and Y. Yuan, "Ccl: Cross-modal correlation learning with multigrained fusion by hierarchical network," *IEEE Transactions on Multimedia*, vol. 20, no. 2, pp. 405–420, 2017.
- [23] J. Qi and Y. Peng, "Cross-modal bidirectional translation via reinforcement learning," in *IJCAI*, 2018, pp. 2630–2636.

- [24] J. Gu, J. Cai, S. R. Joty, L. Niu, and G. Wang, "Look, imagine and match: Improving textual-visual cross-modal retrieval with generative models," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7181–7189.
- [25] W. Gu, X. Gu, J. Gu, B. Li, Z. Xiong, and W. Wang, "Adversary guided asymmetric hashing for cross-modal retrieval," in *Proceedings of the 2019 on international conference on multimedia retrieval*, 2019, pp. 159–167.
- [26] L. He, X. Xu, H. Lu, Y. Yang, F. Shen, and H. T. Shen, "Unsupervised cross-modal retrieval through adversarial learning," in *2017 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2017, pp. 1153–1158.
- [27] B. Wang, Y. Yang, X. Xu, A. Hanjalic, and H. T. Shen, "Adversarial cross-modal retrieval," in *Proceedings of the 25th ACM international conference on Multimedia*, 2017, pp. 154–162.
- [28] L. Zhen, P. Hu, X. Wang, and D. Peng, "Deep supervised cross-modal retrieval," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 10 394–10 403.
- [29] D. Zeng, Y. Yu, and K. Oyama, "Deep triplet neural networks with cluster-cca for audio-visual cross-modal retrieval," *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, vol. 16, no. 3, pp. 1–23, 2020.
- [30] —, "Audio-visual embedding for cross-modal music video retrieval through supervised deep cca," in *2018 IEEE International Symposium on Multimedia (ISM)*. IEEE, 2018, pp. 143–150.
- [31] Q.-Y. Jiang and W.-J. Li, "Deep cross-modal hashing," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 3232–3240.
- [32] W. Wang, X. Yang, B. C. Ooi, D. Zhang, and Y. Zhuang, "Effective deep learning-based multi-modal retrieval," *The VLDB Journal*, vol. 25, no. 1, pp. 79–101, 2016.
- [33] F. Feng, X. Wang, and R. Li, "Cross-modal retrieval with correspondence autoencoder," in *Proceedings of the 22nd ACM international conference on Multimedia*, 2014, pp. 7–16.
- [34] N. Rasiwasia, D. Mahajan, V. Mahadevan, and G. Aggarwal, "Cluster canonical correlation analysis," in *Artificial intelligence and statistics*. PMLR, 2014, pp. 823–831.
- [35] P. L. Lai and C. Fyfe, "Kernel and nonlinear canonical correlation analysis," *International Journal of Neural Systems*, vol. 10, no. 05, pp. 365–377, 2000.
- [36] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [37] Y. Zhou, Z. Wang, C. Fang, T. Bui, and T. L. Berg, "Visual to sound: Generating natural sound for videos in the wild," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3550–3558.
- [38] 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)*, 2018, pp. 247–263.
- [39] S. Hershey, S. Chaudhuri, D. P. Ellis, J. F. Gemmeke, A. Jansen, R. C. Moore, M. Plakal, D. Platt, R. A. Saurous, B. Seybold *et al.*, "Cnn architectures for large-scale audio classification," in *2017 IEEE international conference on acoustics, speech and signal processing (icassp)*. IEEE, 2017, pp. 131–135.
- [40] L. Van der Maaten and G. Hinton, "Visualizing data using t-sne." *Journal of machine learning research*, vol. 9, no. 11, 2008.
- [41] D. Zhai, H. Chang, S. Shan, X. Chen, and W. Gao, "Multiview metric learning with global consistency and local smoothness," *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 3, no. 3, pp. 1–22, 2012.
- [42] X. Zhai, Y. Peng, and J. Xiao, "Heterogeneous metric learning with joint graph regularization for cross-media retrieval," in *Twenty-seventh AAAI conference on artificial intelligence*, 2013.
- [43] J. Wang, Y. He, C. Kang, S. Xiang, and C. Pan, "Image-text cross-modal retrieval via modality-specific feature learning," in *Proceedings of the 5th ACM on International Conference on Multimedia Retrieval*, 2015, pp. 347–354.
- [44] V. Nair and G. E. Hinton, "Rectified linear units improve restricted boltzmann machines," in *Icml*, 2010.
- [45] Y. Peng, X. Huang, and Y. Zhao, "An overview of cross-media retrieval: Concepts, methodologies, benchmarks, and challenges," *IEEE Transactions on circuits and systems for video technology*, vol. 28, no. 9, pp. 2372–2385, 2017.
- [46] S. Akaho, "A kernel method for canonical correlation analysis," *arXiv preprint cs/0609071*, 2006.
- [47] W. Wang and K. Livescu, "Large-scale approximate kernel canonical correlation analysis," *arXiv preprint arXiv:1511.04773*, 2015.
- [48] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.
- [49] H. Zhang, Y. Zhuang, and F. Wu, "Cross-modal correlation learning for clustering on image-audio dataset," in *Proceedings of the 15th ACM international conference on Multimedia*, 2007, pp. 273–276.
- [50] C. Jin, W. Mao, R. Zhang, Y. Zhang, and X. Xue, "Cross-modal image clustering via canonical correlation analysis," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 29, no. 1, 2015.
- [51] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [52] Z. Wang and A. C. Bovik, "Mean squared error: Love it or leave it? a new look at signal fidelity measures," *IEEE signal processing magazine*, vol. 26, no. 1, pp. 98–117, 2009.

- [53] Y. Zhu, Y. Wu, H. Latapie, Y. Yang, and Y. Yan, “Learning audio-visual correlations from variational cross-modal generation,” in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 4300–4304.
- [54] T. J. Mitchell and M. D. Morris, “The spatial correlation function approach to response surface estimation,” in *Proceedings of the 24th conference on Winter simulation*, 1992, pp. 565–571.
- [55] Y. Wen, K. Zhang, Z. Li, and Y. Qiao, “A discriminative feature learning approach for deep face recognition,” in *European conference on computer vision*. Springer, 2016, pp. 499–515.
- [56] 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)*. IEEE, 2017, pp. 776–780.