# Semantically Consistent Multi-view Representation Learning

Yiyang Zhou[a], Qinghai Zheng[b], Shunshun Bai[a], Jihua Zhu[a,*]

[a]*School of Software Engineering, Xi'an Jiaotong University, Xi'an, 710049, China*
[b]*College of Computer and Data Science, Fuzhou University, Fuzhou, 350108, China*

**Abstract**

In this work, we devote ourselves to the challenging task of Unsupervised Multi-view Representation Learning (UMRL), which requires learning a unified feature representation from multiple views in an unsupervised manner. Existing UMRL methods mainly concentrate on the learning process in the feature space while ignoring the valuable semantic information hidden in different views. To address this issue, we propose a novel Semantically Consistent Multi-view Representation Learning (SCMRL), which makes efforts to excavate underlying multi-view semantic consensus information and utilize the information to guide the unified feature representation learning. Specifically, SCMRL consists of a within-view reconstruction module and a unified feature representation learning module, which are elegantly integrated by the contrastive learning strategy to simultaneously align semantic labels of both view-specific feature representations and the learned unified feature representation. In this way, the consensus information in the semantic space can be effectively exploited to constrain the learning process of unified feature representation. Compared with several state-of-the-art algorithms, extensive experiments demonstrate its superiority. Our code is released on https://github.com/YiyangZhou/SCMRL.

*Keywords:* Multi-view representation learning, Contrastive learning, Semantic consensus information

---

*Corresponding author, email: zhujh@xjtu.edu.cn.

## 1. Introduction

Multi-view data are prevalent in real-world applications, and different views can be collected from diverse sensors or various feature extractors. However, numerous classic and effective algorithms [1, 2, 3] are designed for single-view data and can not be leveraged to multi-view data directly. Compared with traditional single-view data, multi-view data are informative and can provide a more comprehensive description[4, 5, 6, 7]. Thanks to these appealing properties, the research of multi-view learning attracts increasing attention, and one of the challenging branches is Unsupervised Multi-view Representation Learning (UMRL). The purpose of UMRL is to learn a unified representation containing consistent information and complementary information, which is usually obtained by mapping data from different views into a shared low-dimensional space[8, 9, 10]. Therefore, the unified representation learned from multiple views could be easily exploited by on-shelf classic single-view algorithms for various downstream tasks effectively[6]. Obviously, a naive way to achieve the goal of UMRL is feature concatenation, which concatenates different views directly to get a joint feature representation. However, since the specific statistical properties among different views are diverse, the feature concatenation strategy usually leads to negative performance[8, 11, 12].

In recent years, many methods are proposed to address the problem of UMRL. For example, CCA[13] and CCA-based methods[14, 15, 16, 17] map different views into a low-dimensional space based on the canonical correlation analysis. FMRL[18] learns the unified representation by utilizing the Hilbert-Schmidt independence criterion [19] to capture the non-linear correlations of multiple views. $AE^2$-Nets[6] introduces the nested autoencoder networks to learn the unified feature representation. DUA-Nets[4] investigates the information of multiple views by employing uncertainty modeling and learns the noise-free feature representation. Although gratifying progress is made and the promising unified multi-view representation can be learned by these aforementioned methods, they are all focused on fusing the multi-view information in the feature space while neglecting the important information in the semantic space. Since different views describe the same object, it is more likely and reasonable to exploit the consensus in the semantic space rather than simply in the feature space during the fusion process[20, 21]. Compared to pursuing consensus in the feature space, exploring consensus in the semantic space can effectively preserve the diversity information of views

2

Figure 1: Overview of SCMRL, which explores and exploits the consensus in the semantic space to boost the learning performance of the unified feature representation of multi-view data. More specifically, given multi-view data $\{\mathbf{X}^{(i)}\}_{i=1}^{m}$ with $m$ views, the within-view reconstruction modules of different views are used to obtain the view-specific feature representation $\{\mathbf{Z}^{(i)}\}_{i=1}^{m}$, and the unified feature representation learning module dynamically learns the unified feature representation $\mathbf{H}$. Meanwhile, the shared classification network and the contrastive learning are introduced to bridge these two modules so that semantic consistency information is captured by $\mathbf{H}$. On the right side, t-SNE is used to visualize the learned $\mathbf{H}$ during the learning process (visualization results are based on the BDGP dataset).

during the learning process of UMRL.

As we discussed above, it is promising to seek the multi-view consensus information in the semantic space to guide the learning process of UMRL. To this end, we propose a novel method, named Semantically Consistent Multi-view Representation Learning (SCMRL), and the framework of SCMRL is depicted in Figure 1. Specifically, SCMRL has two basic modules, namely within-view reconstruction and unified feature representation learning, they are novelty integrated by the exploration of multi-view consensus information in the semantic space.

Different from most existing methods, such as AE$^2$-Nets [6] and CUMRL [22], which mainly consider multi-view information in the feature space, the proposed method employs the multi-view consensus information in the semantic space. As shown in Figure 1, we introduce the shared classification

network to obtain the pseudo labels of both view-specific feature representations $\{\mathbf{Z}^{(i)}\}_{i=1}^m$ and the desired unified feature representation $\mathbf{H}$. The shared classification network acts as an intermediary between within-view reconstruction and unified feature representation learning, which excavates valuable semantic information in the semantic space. Since different views and the learned unified representation describe the same object, their pseudo-labels tend to be consistent. To achieve this goal, we introduce contrastive learning to explore the multi-view consensus information in the semantic space and align the semantic labels of both view-specific feature representations and the learned unified feature representation simultaneously. Based on the aforementioned learning process, the unified representation $\mathbf{H}$ can effectively integrate the information of multiple views in both feature space and semantic space.

The main contributions of the proposed method can be summarized as follows:

• We introduce a novel Semantically Consistent Multi-view Representation Learning (SCMRL), which learns the unified feature representation under the guidance of the consensus semantic information from multi-view data.

• The contrastive learning strategy is specially designed to bridge the within-view reconstruction and unified feature representation learning in SCMRL. The consensus information in the semantic space can be fully exploited to constrain the learning process of the unified feature representation.

• Extensive experiments are conducted, and experimental results verify the effectiveness of the proposed SCMRL compared with several state-of-the-art UMRL methods.

## 2. Related Work

In this section, we will review recent works of UMRL. Since contrastive learning is also leveraged in our method, we will briefly introduce related multi-view learning works based on contrastive learning as well.

### 2.1. UMRL

The goal of UMRL is to learn a promising representation of multi-view data without supervision. The learned unified feature representation from multiple views can be straightforwardly leveraged by off-the-shelf classic and

effective algorithms for downstream tasks, such as classification tasks, clustering tasks, and recognition tasks [6, 17, 4, 23]. Based on the advantage and the effectiveness of multi-view data, UMRL has attracted increasing attention, and many methods have been proposed in recent years [21, 24, 25].

The representative methods are CCA-based[14, 15, 16, 17] methods, which aim to maximize the canonical correlation among views in the low-dimensional space. Taking a multi-view dataset with two views for example here, CCA-based methods have the following basic formula:

$$\min_{\beta_1,\beta_2} -\boldsymbol{I}(\boldsymbol{f}_1(\mathbf{X}^{(1)};\beta_1), \boldsymbol{f}_2(\mathbf{X}^{(2)};\beta_2)) + \lambda\boldsymbol{g}(\beta_1, \beta_2), \qquad (1)$$

where $\boldsymbol{f}_1(\cdot;\beta_1)$ and $\boldsymbol{f}_2(\cdot;\beta_2)$ are two embedding strategies with parameters $\beta_1$ and $\beta_2$. $\boldsymbol{I}(\cdot)$ and $\boldsymbol{g}(\cdot)$ indicate the canonical correlation function and the regularization term respectively. For example, DCCA [15] utilizes deep neural networks for the reconstruction processes of different views. To make the learning process more reliable, DCCAE [17] considers the bottleneck representations by employing autoencoders to minimize the reconstruction loss.

In addition to these CCA-based methods, there are some other UMRL methods that have been proposed in recent years. For example, CMRL[23] fuses the low-dimensional embedding representations and imposes the low-rank tensor constraints on the subspace representations of different views to learn a unified feature representation with comprehensive information. AE$^2$-Nets[6] is a learning framework consisting of nested autoencoders, which is designed to achieve the compact unified multi-view feature representation by balancing the complementarity and consistency properties among views. MvLNet [25] can learn the unified multi-view spectral representation, and uses the Cholesky decomposition during the learning process by introducing the orthogonal constraint and reformulation strategy. DCP[26] is based on information theory, which maximizes mutual information of different views based on contrastive learning to achieve the goal of UMRL.

## 2.2. Contrastive Learning

Contrastive learning[27, 28] is an effective representation learning method, which maximizes the similarity of positive pairs and minimizes the similarity of negative pairs in latent space. This learning paradigm performs well in computer vision[29, 30]. It is also widely used in multi-view clustering in recent years[21, 31, 32, 33]. SURE[34] has designed a contrastive learning

loss specially used for incomplete multi-view clustering, which uses the available pairs as positive pairs and randomly selects some cross-view samples as negative pairs. It effectively and robustly solves the partial view-unaligned problem (PVP) and partial sample-missing problem (PSP) in multi-view clustering. MFLVC[20] optimizes different goals in different feature spaces through contrastive learning and solves the conflict between view reconstruction and consistency goals, which effectively learns common semantics and avoids the influence of meaningless view-specific information.

## 3. Proposed Approach

Given a multi-view dataset $\{\mathbf{X}^{(i)}\}_{i=1}^{m}$ with $m$ view, the $i$-th view in the original feature space is denoted by $\mathbf{X}^{(i)} \in \mathbf{R}^{N \times D_i}$, where $N$ denotes the number of samples and $D_i$ represents the dimension of the feature space. To learn the promising unified multi-view feature representation $\mathbf{H} \in \mathbf{R}^{N \times D_H}$, we propose a novel semantically consistent multi-view representation learning (SCMRL), the framework of which is shown in Figure 1.

### 3.1. SCMRL

Overall, the loss function of SCMRL can be formulated as follows:

$$\mathcal{L} = \mathcal{L}_{Rec} + \lambda_1 \mathcal{L}_{Deg} + \lambda_2 \mathcal{L}_{Sem}, \tag{2}$$

where $\mathcal{L}_{Rec}$ is the loss of within-view reconstruction, $\mathcal{L}_{Deg}$ indicates the loss of the degradation learning employed in unified feature representation learning, and $\mathcal{L}_{Sem}$ denotes the loss of semantic consistency. Regarding $\lambda_1$ and $\lambda_2$, they are two trade-off parameters.

### 3.1.1. Within-view Reconstruction

Generally, the original multi-view data contains a lot of redundant and noisy information, which will have a negative impact on downstream tasks. To simultaneously deal with multi-view data conveniently and learn a reliable representation of each view, we utilize deep autoencoders to obtain the view-specific feature representation of each view. Specifically, $\boldsymbol{E}_i(\cdot; \theta_i)$ and $\boldsymbol{D}_i(\cdot; \phi_i)$ represent the encoder and decoder of the $i$-th view respectively, where $\theta_i$ and $\phi_i$ denote the corresponding parameters.

Subsequently, each view can be encoded into a low-dimensional feature as follows:

$$\mathbf{Z}_j^{(i)} = \boldsymbol{E}_i(\mathbf{X}_j^{(i)}; \theta_i), \tag{3}$$

where $\mathbf{X}_j^{(i)}$ is the $j$-th sample of $\mathbf{X}^{(i)}$, $\mathbf{Z}_j^{(i)} \in \mathbf{R}^{D_Z}$ denotes the embedded feature in the $D_Z$-dimensional feature space. Then we input this low-dimensional feature into the decoder for reconstruction:

$$\hat{\mathbf{X}}_j^{(i)} = \boldsymbol{D}_i(\mathbf{Z}_j^{(i)}; \phi_i), \tag{4}$$

where $\hat{\mathbf{X}}_j^{(i)}$ is the reconstructed representation. Therefore, we can get the following reconstruction loss $\mathcal{L}_{Rec}$:

$$\mathcal{L}_{Rec} = \sum_{i=1}^{m} \sum_{j=1}^{N} ||\mathbf{X}_j^{(i)} - \boldsymbol{D}_i(\boldsymbol{E}_i(\mathbf{X}_j^{(i)}; \theta_i); \phi_i)||_2^2. \tag{5}$$

By minimizing the reconstruction loss $\mathcal{L}_{Rec}$, we can transform the input $\mathbf{X}^{(i)}$ into the representation $\mathbf{Z}^{(i)}$.

### 3.1.2. Unified Feature Representation Learning

Based on the within-view reconstruction, we can obtain the low-dimensional view-specific representations $\{\mathbf{Z}^{(i)}\}_{i=1}^{m}$ of different views. Since the desired unified feature representation should contain comprehensive information of multi-views, it is reasonable that different views require to be reconstructed into a unified feature representation. Considering the way of directly adding the low-dimensional view-specific representation of multiple views, it neglects the different importance and the diverse specific statistical properties of multiple views, and usually leads to poor performance. To effectively learn the unified feature representation $\mathbf{H}$, the degradation learning strategy is adopted in the proposed method. Specifically, we introduce the degradation learning process of the $i$-th view:

$$\mathcal{L}_{Deg} = \sum_{i=1}^{m} \sum_{j=1}^{N} ||\mathbf{Z}_j^{(i)} - \boldsymbol{G}_i(\mathbf{H}_j; \delta_i)||_2^2, \tag{6}$$

where $\boldsymbol{G}_i(\cdot; \delta_i)$ is the fully connected degradation neural network with parameter $\delta_i$ and $\mathbf{H}$ can be updated during the learning process. Based on Eq. (6), the unified feature representation learning can dynamically balance the importance of multiple views and integrate the low-dimensional view-specific representation.

To initialize the unified feature representation, the following strategy is employed:

$$\mathbf{H}_j = \frac{1}{m} \sum_{i=1}^{m} (\mathbf{Z}_j^{(i)}), \tag{7}$$

which can ensure that $\mathbf{H}$ contains comprehensive information of multi-views.

### 3.1.3. Contrastive learning of semantic consistency

As we discussed above, we can observe that both the learning process of within-view reconstruction and unified feature representation learning exploit the multi-view information in the feature space. To effectively excavate the variable semantic consensus information in semantic space, the contrastive learning of semantic consistency is introduced here. According to the fact that multiple views and the unified representation describe the same objective, we introduce a shared classification network, termed $\boldsymbol{Classifer}(\cdot;\varphi)$ with the parameter $\varphi$. Naturally, we constrain that $\{\mathbf{Z}^{(i)}\}_{i=1}^{m}$ and $\mathbf{H}$ should have similar pseudo labels. By utilizing $\boldsymbol{Classifer}(\cdot;\varphi)$, we map the representations, including $\{\mathbf{Z}^{(i)}\}_{i=1}^{m}$ and $\mathbf{H}$, to the semantic space with dimension size of $k$, where $k$ is the number of categories of the multi-view dataset.

In other words, we have the following formula:

$$\{\mathbf{Q}_j^{(i)}, \mathbf{Q}_j^{\mathbf{H}}\} = \boldsymbol{Classifer}(\{\mathbf{Z}_j^{(i)}, \mathbf{H}_j\}; \varphi), \tag{8}$$

in which $\mathbf{Q}_j^{(i)}$ and $\mathbf{Q}_j^{\mathbf{H}}$ denote the pseudo labels of the $\mathbf{Z}_j^{(i)}$ and the unified feature representation $\mathbf{H}_j$, respectively. For convenience, we denote $\mathbf{Q}_j^{(m+1)} = \mathbf{Q}_j^{\mathbf{H}}$. Besides, $\mathbf{Q}_j^{(i)}$ is formulated as:

$$\mathbf{Q}_j^{(i)} = [q_{j1}^{(i)}, q_{j2}^{(i)}...q_{jk}^{(i)}], \tag{9}$$

where $q_{jk}^{(i)}$ is the probability of the $j$-th sample in the $i$-th view belonging to the $k$-th class.

Due to the diverse specific statistic information of multiple views, different views may have confused semantic information in semantic space, which leads to various and confusing results of $\mathbf{Q}_j^{(i)}$ and $\mathbf{Q}_j^{\mathbf{H}}$. Therefore, we introduce contrastive learning [27] to mine the consistent semantic information in the semantic space and obtain consistent categories simultaneously.

Specifically, the sematic column vector $q_{\cdot w}^{(i)}$ have $((m + 1)k - 1)$ vector pairs $\{q_{\cdot w}^{(i)}, q_{\cdot c}^{(j)}\}_{c=1,...,k}^{j=1,...,m+1}$, which contain $m$ positive pairs $\{q_{\cdot w}^{(i)}, q_{\cdot w}^{(j)}\}_{j \neq i}$ and the rest $(k - 1)(m + 1)$ negative pairs. The cosine similarity is utilized to measure the similarity between two semantic column vectors:

$$cos(q_{\cdot c}^{(i)}, q_{\cdot w}^{(j)}) = \frac{q_{\cdot c}^{(i)} \cdot q_{\cdot w}^{(j)}}{||q_{\cdot c}^{(i)}||||q_{\cdot w}^{(j)}||}. \tag{10}$$

---
**Algorithm 1:** Optimization algorithm of SCMRL
---
**Input:** Multi-view dataset $\{\mathbf{X}^{(i)}\}_{i=1}^{m}$; Parameter $\tau$; The number of categories k.

**Output:** Unified Feature Representation $\mathbf{H}$.

1  Initialize $\{\theta_i, \phi_i\}_{i=1}^{m}$ by minimizing Eq. (5);
2  Initialize $\mathbf{H}$ by Eq. (7);
3  **while** not converged **do**
4      Obtain the view-specific representation $\{\mathbf{Z}^{(i)}\}_{i=1}^{m}$ through Eq. (3);
5      Obtain the label distribution of each view and the unified feature representation by Eq. (8);
6      Update $\{\theta_i, \phi_i, \delta_i\}_{i=1}^{m}$, $\mathbf{H}$ and $\varphi$ with Eq. (2);
7  **end**
---

Then, contrastive loss $l_{two}(i, j)$ between semantic column vectors $q_{\cdot c}^{(i)}$ and $q_{\cdot c}^{(j)}$ is defined as:

$$l_{two}(i, j) = -\sum_{c=1}^{k} log(f(i, j, c)), \tag{11}$$

where

$$f(i, j, c) = \frac{e^{cos(q_{\cdot c}^{(i)}, q_{\cdot c}^{(j)})/\tau}}{\left(\sum_{w=1}^{k}\left(e^{cos(q_{\cdot c}^{(i)}, q_{\cdot w}^{(i)})/\tau} + e^{cos(q_{\cdot c}^{(i)}, q_{\cdot w}^{(j)})/\tau}\right) - e^{\frac{1}{\tau}}\right)} \tag{12}$$

and $\tau$ denotes the temperature parameter. Further, it is convenient to define the complete comparative learning loss of semantic consistency $L_{Sem}$ as:

$$L_{Sem} = l_{Sum} + l_{Reg}, \tag{13}$$

where $l_{Sum}$ denotes the contrastive loss for the whole dataset and $l_{Reg}$ indicates the regularization phase. Specifically, the item $l_{Sum}$ is defined as:

$$l_{Sum} = \frac{1}{2}\sum_{i=1}^{m+1}\sum_{j=1, j\neq i}^{m+1}\frac{l_{two}(i, j)}{k}. \tag{14}$$

And the item $l_{Reg}$ is formulated as:

$$l_{Reg} = \sum_{i=1}^{m+1}\sum_{c=1}^{k}\left(\frac{1}{N}\sum_{j=1}^{N}q_{jc}^{(i)} \log \frac{1}{N}\sum_{j=1}^{N}q_{jc}^{(i)}\right), \tag{15}$$

9

Table 1: Details of the used datasets.

| Dataset | #Sample | #Cluster | #View | #Dimensionality of features |
|---|---|---|---|---|
| MNIST-USPS | 5000 | 10 | 2 | {1*28*28, 1*28*28} |
| BDGP | 2500 | 5 | 2 | {1750, 79} |
| Fashion | 10000 | 10 | 3 | {1*28*28, 1*28*28, 1*28*28} |
| CCV | 6773 | 20 | 3 | {5000, 5000, 4000} |
| Caltech-2V | 1400 | 7 | 2 | {40, 254} |
| Caltech-3V | 1400 | 7 | 3 | {40, 254, 928} |
| Caltech-4V | 1400 | 7 | 4 | {40, 254, 928, 512} |
| Caltech-5V | 1400 | 7 | 5 | {40, 254, 928, 512, 1984} |

which can avoid grouping all samples into the same cluster.

For clarification, the optimization procedure of SCMRL is summarized in Algorithm 1.

## 4. Experiments

To verify the effectiveness of our method, extensive experiments are conducted in this section. Specifically, two basic tasks, namely clustering and classification, are used to evaluate the performance of the learned unified multi-view feature representation. Furthermore, detailed discussions of our method are provided as well.

### 4.1. Experiments Setup

#### 4.1.1. Datasets

We use the following benchmark datasets, the detail of these datasets are shown in Table 1:

**1) MNIST-USPS**[35]: It is a two-view dataset that contains 5000 handwritten digital image samples from numbers 0 to 9.

**2) BDGP**[36]: It is a two-view dataset containing 2500 images of drosophila embryos belonging to 5 categories.

**3) Fashion**[37]: It has 10000 images collected from 10 categories about fashion products and has three views.

**4) CCV**[38]: It contains 6773 internet videos samples belonging to 20 classes. It has three views, such as STIP, SIFT, and MFCC.

**5) Caltech**[39]: It is collected from 1400 images, which belong to 7 categories and have five views. Four sub-datasets, namely Caltech-2V, Caltech-3V, Caltech-4V, and Caltech-5V, with different numbers of views are built for evaluation here. Specifically, Caltech-2V uses WM and CENTRISTT; Caltech-3V uses WM, CENTRIST, and LBP; Caltech-4V uses WM, CENTRIST, LBP, and GIST; Caltech-5V uses WM, CENTRIST, LBP, GIST, and HOG.

### 4.1.2. Comparison methods

The following state-of-the-arts are used for comparison:

**1) DCCA**[15]: It is a classic CCA-based method, which uses depth neural networks to obtain the nonlinear mapping with the maximum view linear correlation.

**2) DCCAE**[17]: It is also a CCA-based method. Different from DCCA, it uses autoencoders to obtain the low dimensional projection of the original data and maximizes the view correlation between the learned representations.

**3) LMSC**[40]: It learns the latent data representation by mapping different views into the common space and employing the low-rank subspace constraint.

**4) AE$^2$-Nets**[6]: Nested autoencoders are used to learn the compact unified representation by balancing the complementarity and consistency among multi-views.

**5) DUA-Nets**[4]: It presents the dynamic uncertainty-aware networks for UMRL. By estimating and leveraging the uncertainty of data, it achieves the noise-free multi-view feature representation.

**6) CUMRL**[22]: It considers the low-rank tensor constraint to excavate the high-order view correlations of multi-view data in the feature space, and introduces the collaborative learning strategy for UMRL.

**7) DCP**[26]: It learns the unified multi-view representation by maximizing the mutual information of different views via contrastive learning in the feature space.

**8) CMRL**[23]: It achieves the unified multi-view representation with comprehensive information by introducing the orthogonal mapping strategy and imposing the low-rank tensor constraint on the subspace representations.

### 4.1.3. Evaluation metrics

Two basic tasks, i.e., clustering and classification, are employed in this section. Since $k$-means and $k$-nearest neighbor ($k$NN) are simple and can

(a) Epoch 10        (b) Epoch 20

(c) Epoch 30      (d) Epoch 40      (e) Epoch 50

Figure 2: The t-SNE visualization results of the learned unified feature representation in Epoch 10, 20, 30, 40, and 50 of the learning process. Experimental results on the MNIST-USPS dataset are shown here.

effectively measure the quality of the learned unified representation, we adopt the classic $k$-means algorithm for the clustering task and $k$NN algorithm for the classification task here[6, 12, 4].

**1) Clustering task**: Three metrics are utilized to evaluate the clustering quality, namely ACC, Normalized Mutual Information (NMI), and Fscore. For each dataset, 10 trials are conducted for all experiments to eliminate the randomness and make the experimental results more reliable.

**2) Classification task**: We utilize ACC as the metric to evaluate the classification performance. Specifically, we set k=5 for $k$NN in all experiments and 30 trials are conducted for each experiment. We divide samples of the learned representation into training and testing sets with different proportions, and the ratios (#Train/#Test) are set to 8/2 (80%/20%), 5/5 (50%/50%), 2/8 (20%/80%).

Table 2: Clustering results of the learned unified multi-view representation on BDGP, MNIST-USPS, Fashion and CCV.

| Datasets | BDGP | | | MNIST-USPS | | | Fashion | | | CCV | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Metrics | ACC | NMI | Fscore | ACC | NMI | Fscore | ACC | NMI | Fscore | ACC | NMI | Fscore |
| DCCA | 67.72 | 56.45 | 63.05 | 41.55 | 35.13 | 31.52 | 67.38 | 77.51 | 72.66 | 21.09 | 19.14 | 13.26 |
| DCCAE | 69.03 | 58.12 | 63.31 | 67.94 | 53.62 | 51.60 | 67.74 | 77.26 | 72.32 | 20.42 | 18.56 | 13.00 |
| LMSC | 52.39 | 43.15 | 55.63 | 37.25 | 43.33 | 40.42 | 43.90 | 41.68 | 36.71 | 13.49 | 8.660 | 11.76 |
| AE$^2$-Nets | 55.24 | 40.57 | 50.13 | 62.59 | 62.31 | 56.71 | 72.94 | 76.27 | 71.57 | 9.430 | 2.810 | 7.590 |
| DUA-Nets | 60.28 | 40.61 | 53.91 | 75.10 | 68.94 | 66.01 | 77.21 | 76.08 | 72.72 | 16.11 | 11.80 | 11.64 |
| CUMRL | 62.93 | 48.80 | 57.44 | 58.64 | 56.93 | 50.96 | 67.00 | 66.77 | 61.26 | 10.41 | 5.190 | 7.480 |
| DCP | 43.77 | 38.50 | 53.49 | 89.10 | 94.13 | 92.87 | 75.68 | 86.19 | 82.23 | 14.23 | 11.48 | 10.69 |
| CMRL | 78.92 | 67.20 | 71.68 | 91.61 | 85.60 | 85.98 | 76.81 | 80.29 | 74.77 | 24.31 | 18.11 | 13.64 |
| Ours | **98.00** | **94.69** | **96.18** | **99.56** | **98.72** | **99.12** | **98.90** | **97.32** | **97.84** | **26.82** | **27.55** | **21.39** |

Table 3: Clustering results of the learned unified multi-view representation on Caltech-$n$V, $n$ is selected from $\{2, 3, 4, 5\}$.

| Datasets | Caltech-2V | | | Caltech-3V | | | Caltech-4V | | | Caltech-5V | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Metrics | ACC | NMI | Fscore | ACC | NMI | Fscore | ACC | NMI | Fscore | ACC | NMI | Fscore |
| DCCA | 39.36 | 30.69 | 37.02 | 46.71 | 35.29 | 37.78 | 54.97 | 33.44 | 38.53 | 62.73 | 42.97 | 48.13 |
| DCCAE | 43.87 | 31.17 | 36.23 | 59.49 | 44.53 | 47.15 | 50.29 | 32.01 | 37.78 | 63.70 | 45.58 | 49.62 |
| LMSC | 42.64 | 32.35 | 35.47 | 26.44 | 7.580 | 20.32 | 38.08 | 28.82 | 34.43 | 66.15 | 53.38 | 56.83 |
| AE$^2$-Nets | 46.14 | 32.01 | 35.60 | 51.48 | 41.08 | 42.89 | 48.01 | 38.89 | 41.93 | 67.67 | 58.13 | 57.96 |
| DUA-Nets | 39.45 | 22.12 | 30.88 | 43.63 | 29.37 | 36.96 | 46.31 | 34.54 | 41.03 | 56.89 | 44.37 | 48.46 |
| CUMRL | 48.74 | 41.32 | 42.36 | 56.59 | 48.86 | 49.29 | 67.76 | 57.66 | 57.81 | 88.55 | 79.02 | 81.12 |
| DCP | 42.64 | 32.35 | 35.47 | 51.60 | 48.35 | 52.16 | 53.37 | 53.55 | 55.54 | 54.04 | 54.05 | 56.16 |
| CMRL | 55.03 | 40.33 | 42.95 | 59.27 | 44.67 | 46.72 | 69.30 | 56.24 | 57.75 | 68.76 | 56.37 | 57.61 |
| Ours | **61.29** | **48.29** | **48.75** | **78.21** | **70.18** | **70.30** | **87.00** | **79.83** | **79.22** | **89.00** | **80.23** | **81.22** |

### 4.1.4. Implementation details

For all datasets, the ReLU[41] activation function is used to implement autoencoders in SCMRL. Adam optimizer[42] is employed for optimization. Our method is implemented by PyTorch[43] on one NVIDIA Geforce GTX 2080ti GPU with 11GB memory. For comparison methods, we leverage the codes released by their corresponding authors and use the recommended settings in their original works.

### 4.2. Visualization results

To vividly reveal the structure of the learned unified representation, we visualize **H** achieved in the Epoch 10, 20, 30, 40, and 50 of the SCMRL learning process based on the t-SNE [44]. Taking the MNIST-USPS dataset for example here, the visualization results are shown in Figure 2. It can be observed that the unified feature representation with a promising structure can be achieved by our method.

13

Table 4: Classification results of the learned unified multi-view representation on BDGP, MNIST-USPS, Fashion and CCV.

| Datasets | BDGP | | | MNIST-USPS | | | Fashion | | | CCV | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| #Train/#Test | 8/2 | 5/5 | 2/8 | 8/2 | 5/5 | 2/8 | 8/2 | 5/5 | 2/8 | 8/2 | 5/5 | 2/8 |
| DCCA | 96.52 | 96.44 | 95.25 | 78.48 | 76.46 | 72.07 | 86.93 | 86.47 | 85.56 | 28.49 | 26.64 | 24.01 |
| DCCAE | 97.73 | 97.40 | 96.62 | 83.45 | 81.74 | 79.44 | 86.63 | 86.11 | 85.24 | 29.71 | 27.97 | 24.87 |
| LMSC | 98.16 | 97.60 | 96.22 | 94.66 | 93.34 | 89.55 | 86.12 | 84.80 | 81.27 | 33.77 | 31.17 | 26.02 |
| AE$^2$-Nets | 91.65 | 91.19 | 87.56 | 96.43 | 95.62 | 93.30 | 91.83 | 90.95 | 89.10 | 6.380 | 6.280 | 6.070 |
| DUA-Nets | 93.01 | 91.52 | 85.42 | 91.79 | 90.18 | 85.64 | 87.60 | 86.53 | 81.47 | 30.43 | 31.04 | 24.86 |
| CUMRL | 94.73 | 93.57 | 90.91 | 94.69 | 93.99 | 91.16 | 81.49 | 80.61 | 79.26 | 6.75 | 6.67 | 6.49 |
| DCP | 96.38 | 96.04 | 93.94 | <u>98.49</u> | <u>98.33</u> | <u>97.95</u> | 93.60 | 93.40 | <u>92.44</u> | 20.25 | 18.45 | 15.64 |
| CMRL | <u>98.82</u> | <u>98.57</u> | <u>98.24</u> | 97.86 | 97.45 | 96.28 | <u>93.85</u> | <u>93.48</u> | 92.36 | <u>36.03</u> | <u>34.18</u> | <u>30.63</u> |
| Ours | **98.92** | **98.74** | **98.54** | **99.63** | **99.57** | **99.57** | **99.11** | **99.08** | **99.04** | **39.41** | **37.85** | **35.09** |

## 4.3. Experimental Results

We discuss the experimental results of the clustering task and the classification task. Overall speaking, the proposed method can achieve the best performance in most cases.

### 4.3.1. Experimental results of the clustering task

The $k$-means clustering results on all datasets are shown in Table 2 and 3. We can observe that SCMRL achieves the best performance on all datasets in all metrics and considerable progress can be made for all metrics. For example, on the BDGP dataset, our method respectively achieves an improvement of around 19.08%, 27.49%, and 24.50% compared with the second-best results in metrics of ACC, NMI, and Fscore. Compared with the second-best results, 17.70%, 22.17%, and 21.41% improvements can be obtained on the Caltech-4V in metrics of ACC, NMI, and Fscore, respectively. The main reason is that the proposed SCMRL explores the underlying consistent information of both the learned unified representation and multiple views in semantic space. Compared with other methods, our method can explore the consistent information in semantic space. Consequently, the multi-view diversity information in the feature space can be preserved and utilized during the learning process of the unified feature representation. In Table 3, we can find that the performance of SCMRL steadily increases with the increase of views on Caltech-$n$V, which indicates that the rich information contained in multiple views can be effectively excavated and integrated into the learned unified representation by our method.

14

Table 5: Classification results of the learned unified multi-view representation on Caltech-$n$V, $n$ is selected from $\{2, 3, 4, 5\}$.

| Datasets | Caltech-2V | | | Caltech-3V | | | Caltech-4V | | | Caltech-5V | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| #Train/#Test | 8/2 | 5/5 | 2/8 | 8/2 | 5/5 | 2/8 | 8/2 | 5/5 | 2/8 | 8/2 | 5/5 | 2/8 |
| DCCA | 68.24 | 65.41 | 60.61 | 72.81 | 71.14 | 66.81 | 65.80 | 64.26 | 61.04 | 62.73 | 42.97 | 48.13 |
| DCCAE | 69.82 | 67.12 | 62.92 | 74.46 | 73.01 | 68.23 | 66.80 | 64.75 | 60.19 | 63.70 | 45.58 | 49.62 |
| LMSC | 72.11 | 70.22 | 63.29 | 72.25 | 69.44 | 60.55 | 76.71 | 73.81 | 67.01 | 90.89 | 89.87 | 86.31 |
| AE$^2$-Nets | 69.80 | 66.78 | 61.09 | 82.93 | 80.89 | 76.27 | 84.83 | 84.19 | 79.34 | 91.32 | 90.58 | 88.04 |
| DUA-Nets | 62.73 | 59.40 | 54.37 | 69.58 | 69.28 | 63.05 | 75.88 | 74.29 | 66.69 | 81.04 | 80.29 | 77.51 |
| CUMRL | <u>80.21</u> | <u>77.64</u> | <u>71.89</u> | **87.18** | **85.83** | <u>81.62</u> | <u>90.31</u> | <u>89.70</u> | <u>86.59</u> | <u>93.48</u> | <u>93.29</u> | <u>91.83</u> |
| DCP | 71.36 | 70.68 | 67.23 | 80.75 | 79.41 | 77.15 | 82.74 | 81.33 | 78.04 | 87.81 | 86.21 | 83.46 |
| CMRL | **80.81** | **79.05** | **74.10** | 84.74 | 82.99 | 79.46 | 90.27 | 89.21 | 86.23 | 91.01 | 89.90 | 86.71 |
| Ours | 73.80 | 71.90 | 68.27 | <u>84.82</u> | <u>83.30</u> | **81.65** | **91.07** | **90.46** | **88.78** | **94.15** | **93.95** | **91.93** |

## 4.3.2. Experimental results of the classification task

The classification results are shown in Table 4 and Table 5. In general, the proposed method can obtain promising results for all datasets and achieve the best results for most cases. For example, on the Fashion dataset, our method can achieve an improvement of around 7% with respect to the metric of ACC, compared with the second-best results. Although CMRL and CUMRL achieve slightly better classification results on Caltech-2V and Caltech-3V, the results of our method are also competitive. Furthermore, with the increase of view on Caltech-$n$V, the performance of SCMRL can improve more significantly than CMRL and CUMRL. With the decrease of #Train/#Test, it can be observed that our method has the slowest performance decline, which also indicates the effectiveness of the semantic consistent information exploration in our method.

Table 6: Ablation study on Fashion dataset. "✓" indicates the used component, $Clu$ and $Cla$ denote the $Clu$stering task and the $Cla$ssification task respectively.

| Dataset | $\mathcal{L}_{Sem}$ | $\mathcal{L}_{Rec}$ | $Clu$ | | $Cla$ |
|---|---|---|---|---|---|
| | | | NMI | Fscore | ACC |
| | ✓ | ✓ | **98.72** | **99.12** | **99.63** |
| Fashion | | ✓ | 81.69 | 76.19 | 93.57 |
| | ✓ | | 96.74 | 97.34 | 98.49 |

(a) ACC in classification task



(b) NMI in clustering task



(c) Sensitivity of model to parameter $\tau$ in classification task



(d) Sensitivity of model to parameter $\tau$ in clustering task

Figure 3: Parameter sensitivity analysis on MNIST-USPS dataset.

## 4.4. Model Analysis

### 4.4.1. Ablation Studies

We conduct ablation studies here. It is clear that the module of unified feature representation learning is essential since UMRL aims to learn the unified multi-view representation. Consequently, we discuss the learning process of our method with and without $\mathcal{L}_{Sem}$ and $\mathcal{L}_{Rec}$. We take the experiments on the Fashion dataset for example. The clustering results (in metrics of NMI and Fscore) and the classification results (in the metric of ACC with #Train/#Test = 80%/20%) are reported in Table 6. According to the experimental results, we conclude: 1) Both the employment of $\mathcal{L}_{Sem}$ and $\mathcal{L}_{Rec}$ effectively improve the learned unified representation; 2) Compared with the employment of $\mathcal{L}_{Rec}$, the employment of $\mathcal{L}_{Sem}$ can significantly boost the

Table 7: Performance comparison between single view and multi-view.

| Dataset | MNIST-USPS | | BDGP | |
|---|---|---|---|---|
| Metrics | ACC | NMI | ACC | NMI |
| $k$-means (view 1) | 76.78 | 72.33 | 45.00 | 26.49 |
| $k$-means (view 2) | 58.10 | 52.13 | 57.04 | 45.96 |
| DEC (view 1) | 73.10 | 71.49 | 46.28 | 29.96 |
| DEC (view 2) | 56.12 | 61.08 | 94.78 | 86.62 |
| SCMRL_concat | 82.36 | 76.98 | 71.00 | 62.03 |
| SCMRL_average | 89.48 | 91.62 | 81.72 | 67.22 |
| SCMRL | **99.56** | **98.72** | **98.00** | **94.69** |

performance of the learned unified representation.

In addition, we further discussed how much performance improvement of multi-view data compared with single-view data on downstream tasks and how different ways of integrating multi-view information affect the performance of downstream tasks. Taking the clustering task as an example, we compare SCMRL and its different variants with the classical single-view clustering algorithm [3] on BDGP and MNIST-USPS, and the results are shown in Table 7. In which, we use SCMRL_concat to directly splice the low-dimensional representations $\{\mathbf{Z}^{(i)}\}_{i=1}^{m}$ of different views for clustering, while SCMRL_average is to average $\{\mathbf{Z}^{(i)}\}_{i=1}^{m}$ of different views for clustering.

Convincingly, these aforementioned observations of ablation studies indicate that introducing excavation of multi-view semantic consistent information in the semantic space plays a vital role in our method for the learning process of multi-view unified feature representation. Multi-view data has more descriptive information than single-view data, which is helpful for downstream tasks. At the same time, in order to better integrate the data of multiple views, SCMRL's unique fusion method is more effective than simple splicing and averaging.

*4.4.2.* **Parameter sensitivity analysis**

To explore the sensitivity of SCMRL to hyper-parameters, we take different values of $\lambda_1$ and $\lambda_2$ in Eq. (2) on MNIST-USPS dataset and explore their influence of the clustering task (in the metric of NMI) and the classification task (in the metric of ACC). In order to make the results more reliable, we run the clustering task and the classification task 10 times and 30 times respectively to average. The results are shown in Figure 3(a) and Figure 3(b),

which indicates that our model is insensitive to $\lambda_1$ and $\lambda_2$. Actually, $\lambda_1$ and $\lambda_2$ are set to 1 for all datasets in Table 2-5.

As for the hyper-parameter $\tau$ in Eq. (11), results of the MNIST-USPS datasets are shown in Figure 3(c) and Figure 3(d). It can also be found that SCMRL is also robust to $\tau$. Actually, we fix parameter $\tau = 0.5$ for all datasets in Table 2-5.



Figure 4: Convergence curve on MNIST-USPS dataset.

### *4.4.3.* **Convergence analysis**

To show the convergence properties of SCMRL, we take the experiment on MNIST-USPS dataset for example and display the convergence curve in Figure 4. It can be observed that the loss value drops rapidly in the first 20 epochs and has promising convergence properties. For other datasets, the similar convergence properties can be achieved as well.

## 5. Conclusion

In this paper, we introduce a novel Semantically Consistent Multi-view Representation Learning (SCMRL), which effectively excavates and exploits

the consistent information of multiple views in the semantic space to learn unified multi-view feature representation with promising structure. By introducing the contrastive learning of semantic consistency, SCMRL incorporates the within-view reconstruction with the unified feature representation learning and explores the valuable consensus information in semantic space to guide the learning process. Experimental results conducted on several benchmark datasets verify the effectiveness of SCMRL over other state-of-the-art methods.

## References

[1] J. Shi, J. Malik, Normalized cuts and image segmentation, IEEE Transactions on pattern analysis and machine intelligence 22 (8) (2000) 888–905.

[2] J. MacQueen, Classification and analysis of multivariate observations, in: 5th Berkeley Symp. Math. Statist. Probability, 1967, pp. 281–297.

[3] J. Xie, R. Girshick, A. Farhadi, Unsupervised deep embedding for clustering analysis, in: International conference on machine learning, PMLR, 2016, pp. 478–487.

[4] Y. Geng, Z. Han, C. Zhang, Q. Hu, Uncertainty-aware multi-view representation learning, in: Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 35, 2021, pp. 7545–7553.

[5] C. Zhang, H. Fu, Q. Hu, X. Cao, Y. Xie, D. Tao, D. Xu, Generalized latent multi-view subspace clustering, IEEE transactions on pattern analysis and machine intelligence 42 (1) (2018) 86–99.

[6] C. Zhang, Y. Liu, H. Fu, Ae2-nets: Autoencoder in autoencoder networks, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2019, pp. 2577–2585.

[7] Q. Zheng, J. Zhu, Z. Li, S. Pang, J. Wang, Y. Li, Feature concatenation multi-view subspace clustering, Neurocomputing 379 (2020) 89–102.

[8] C. Xu, D. Tao, C. Xu, A survey on multi-view learning, arXiv preprint arXiv:1304.5634 (2013).

[9] C. Lu, J. Feng, Y. Chen, W. Liu, Z. Lin, S. Yan, Tensor robust principal component analysis with a new tensor nuclear norm, IEEE transactions on pattern analysis and machine intelligence 42 (4) (2019) 925–938.

[10] L. Qu, M. Liu, D. Cao, L. Nie, Q. Tian, Context-aware multi-view summarization network for image-text matching, in: Proceedings of the 28th ACM International Conference on Multimedia, 2020, pp. 1047–1055.

[11] J. Zhao, X. Xie, X. Xu, S. Sun, Multi-view learning overview: Recent progress and new challenges, Information Fusion 38 (2017) 43–54.

[12] Q. Zheng, J. Zhu, Z. Li, H. Tang, Graph-guided unsupervised multi-view representation learning, IEEE Transactions on Circuits and Systems for Video Technology (2022).

[13] H. Hotelling, Relations between two sets of variates, in: Breakthroughs in statistics, Springer, 1992, pp. 162–190.

[14] S. Akaho, A kernel method for canonical correlation analysis, arXiv preprint cs/0609071 (2006).

[15] G. Andrew, R. Arora, J. Bilmes, K. Livescu, Deep canonical correlation analysis, in: International conference on machine learning, PMLR, 2013, pp. 1247–1255.

[16] J. Chen, G. Wang, G. B. Giannakis, Graph multiview canonical correlation analysis, IEEE Transactions on Signal Processing 67 (11) (2019) 2826–2838.

[17] W. Wang, R. Arora, K. Livescu, J. Bilmes, On deep multi-view representation learning, in: International conference on machine learning, PMLR, 2015, pp. 1083–1092.

[18] R. Li, C. Zhang, Q. Hu, P. Zhu, Z. Wang, Flexible multi-view representation learning for subspace clustering., in: IJCAI, 2019, pp. 2916–2922.

[19] A. Gretton, O. Bousquet, A. Smola, B. Schölkopf, Measuring statistical dependence with hilbert-schmidt norms, in: International conference on algorithmic learning theory, Springer, 2005, pp. 63–77.

[20] J. Xu, H. Tang, Y. Ren, L. Peng, X. Zhu, L. He, Multi-level feature learning for contrastive multi-view clustering, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2022, pp. 16051–16060.

[21] Y. Li, P. Hu, Z. Liu, D. Peng, J. T. Zhou, X. Peng, Contrastive clustering, in: Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 35, 2021, pp. 8547–8555.

[22] Q. Zheng, J. Zhu, Z. Li, Collaborative unsupervised multi-view representation learning, IEEE Transactions on Circuits and Systems for Video Technology (2021).

[23] Q. Zheng, J. Zhu, Z. Li, Z. Tian, C. Li, Comprehensive multi-view representation learning, Information Fusion 89 (2023) 198–209.

[24] J. Liu, X. Liu, Y. Yang, X. Guo, M. Kloft, L. He, Multiview subspace clustering via co-training robust data representation, IEEE Transactions on Neural Networks and Learning Systems (2021).

[25] Z. Huang, J. T. Zhou, H. Zhu, C. Zhang, J. Lv, X. Peng, Deep spectral representation learning from multi-view data, IEEE Transactions on Image Processing 30 (2021) 5352–5362.

[26] Y. Lin, Y. Gou, X. Liu, J. Bai, J. Lv, X. Peng, Dual contrastive prediction for incomplete multi-view representation learning, IEEE Transactions on Pattern Analysis and Machine Intelligence (2022).

[27] T. Chen, S. Kornblith, M. Norouzi, G. Hinton, A simple framework for contrastive learning of visual representations, in: International conference on machine learning, PMLR, 2020, pp. 1597–1607.

[28] E. Xie, J. Ding, W. Wang, X. Zhan, H. Xu, P. Sun, Z. Li, P. Luo, Detco: Unsupervised contrastive learning for object detection, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, 2021, pp. 8392–8401.

[29] C. Niu, H. Shan, G. Wang, Spice: Semantic pseudo-labeling for image clustering, arXiv preprint arXiv:2103.09382 (2021).

[30] W. Van Gansbeke, S. Vandenhende, S. Georgoulis, M. Proesmans, L. Van Gool, Scan: Learning to classify images without labels, in: European conference on computer vision, Springer, 2020, pp. 268–285.

[31] Y. Lin, Y. Gou, Z. Liu, B. Li, J. Lv, X. Peng, Completer: Incomplete multi-view clustering via contrastive prediction, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2021, pp. 11174–11183.

[32] S. Roy, A. Etemad, Self-supervised contrastive learning of multi-view facial expressions, in: Proceedings of the 2021 International Conference on Multimodal Interaction, 2021, pp. 253–257.

[33] F. Lin, B. Bai, K. Bai, Y. Ren, P. Zhao, Z. Xu, Contrastive multi-view hyperbolic hierarchical clustering, arXiv preprint arXiv:2205.02618 (2022).

[34] M. Yang, Y. Li, P. Hu, J. Bai, J. C. Lv, X. Peng, Robust multi-view clustering with incomplete information, IEEE Transactions on Pattern Analysis and Machine Intelligence (2022).

[35] X. Peng, Z. Huang, J. Lv, H. Zhu, J. T. Zhou, Comic: Multi-view clustering without parameter selection, in: International conference on machine learning, PMLR, 2019, pp. 5092–5101.

[36] X. Cai, H. Wang, H. Huang, C. Ding, Joint stage recognition and anatomical annotation of drosophila gene expression patterns, Bioinformatics 28 (12) (2012) i16–i24.

[37] H. Xiao, K. Rasul, R. Vollgraf, Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms, arXiv preprint arXiv:1708.07747 (2017).

[38] Y.-G. Jiang, G. Ye, S.-F. Chang, D. Ellis, A. C. Loui, Consumer video understanding: A benchmark database and an evaluation of human and machine performance, in: Proceedings of the 1st ACM International Conference on Multimedia Retrieval, 2011, pp. 1–8.

[39] L. Fei-Fei, R. Fergus, P. Perona, Learning generative visual models from few training examples: An incremental bayesian approach tested on 101

object categories, in: 2004 conference on computer vision and pattern recognition workshop, IEEE, 2004, pp. 178–178.

[40] C. Zhang, Q. Hu, H. Fu, P. Zhu, X. Cao, Latent multi-view subspace clustering, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 4279–4287.

[41] X. Glorot, A. Bordes, Y. Bengio, Deep sparse rectifier neural networks, in: Proceedings of the fourteenth international conference on artificial intelligence and statistics, JMLR Workshop and Conference Proceedings, 2011, pp. 315–323.

[42] D. P. Kingma, J. Ba, Adam: A method for stochastic optimization, arXiv preprint arXiv:1412.6980 (2014).

[43] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al., Pytorch: An imperative style, high-performance deep learning library, Advances in neural information processing systems 32 (2019).

[44] L. Van der Maaten, G. Hinton, Visualizing data using t-sne., Journal of machine learning research 9 (11) (2008).