# Improving Image Recognition Accuracy Using Multi-Views Spectral Clustering

Saja Hikmat Dawood

Department of Computer, College of Basic Education, Mustansiriyah University, Baghdad 10052, Iraq

Corresponding Author Email: phd202130684@iips.edu.iq

**ABSTRACT**

This paper presents a new way to tremendously improve the picture clustering quality by exploiting multiple "views" of the data. Image grouping is a process of grouping photographs associated with visual characteristics. The most appropriate characteristics and AI architectures for picture clustering have, to date, been very difficult to select although they have a significant impact on the quality of the clustering results. As a solution to the challenge, the so-called Multi-View Clustering (MVC) is proposed. In MVC, multiple AI networks, which are already pre-trained, like convolutional neural networks (CNNs), act as multiple "views" with regards to the same visual information. While drawing from the same data, each of these CNNs captures another perspective by extracting a unique set of image features. This method will attempt to consider various points of view in order to collect different and complementary information about the images. A neural network architecture with multiple inputs is proposed for this many-views problem. Trained end-to-end, this resolves the MVC problem using the features extracted from each of the CNN views as input. Improved pooling performance is a result of end-to-end training that ensures the network has learnt how to aggregate features from multiple views efficiently. Experimental results on several image datasets have proven the usefulness of this strategy. The proposed method is with an end-to-end training strategy, utilizing several jointly pre-trained CNNs as feature extractors, so it outperforms conventional image clustering accuracy. Indeed, state-of-the-art results are produced in the field of image collage. Conclusively, this paper proposes a holistic approach that improves the efficiency of image classification, which is a critical contribution to the literature on image clustering. The approach proposed overcomes the challenge of feature selection and AI architecture for the image clustering through the use of multiple views of spectral ensembles to some pre-trained AI networks and leveraging an end-to-end training approach. The findings show how the efficiency of picture grouping methods is improved through the incorporation of numerous viewpoints.

## 1. INTRODUCTION

Image clustering is an important technology in computer vision, which can help organize and sort a large amount of visual data. Categorizing the pictures based on their visual characteristics can significantly improve the understanding of picture content and patterns. The choice of what image features and AI architectures are best is a common issue in many traditional images clustering techniques, leading to poor clustering quality [1].

Traditional image clustering approaches, such as the K-means method and hierarchical clustering, have been largely depending on handcrafted features, which often lack the ability to capture complex data patterns like images. These methodologies are not suitable for high-dimensional data and they often assume the data to be linearly separable. These assumptions and limitations are known to degrade the effectiveness and scalability of these algorithms significantly. Multi-View Spectral Clustering (MVSC), on the other hand, overcomes these limitations by using different representations of the data. This approach merges diverse views (or representations) of the data, where each one corresponds to a different feature descriptor of the image, to model the similarities between them.

MVSC outperforms traditional image clustering methodologies by incorporating multiple, complementary sets of features to capture content and patterns in the data. It also utilizes a wide spectrum of information from different views (features) rather than just a single, possibly inadequate representation and therefore provides more robust and accurate clustering results [2].

The creation of MVSC has been crucial for mitigating these issues and improving the accuracy of picture categorization. To enhance the AI networks' understanding of pictures, it makes use of multiple "views" of the information stored in pictures. The fundamental idea of MVSC is to use trained AI systems—CNNs in this instance—each as separate "views" on the image set. Each CNN captures different perspectives or representations of the data, and each extracts a different set of characteristics, in order to offer a broad understanding of the visual information contained within a picture. By using different AI networks [2], it can harvest complementary and

nuanced information about the images. In order to effectively use the strengths of the different views that provide different representations, they propose to utilize a multi-input neural architecture. This innovative design addresses the Multi-View Spectral Clustering issue by using the feature set derived from each CNN view as input and supports an end-to-end training approach. The network gain from an end-to- end training strategy because the network learns how to better combine representations from different views, and ultimately - by dissecting the images in different ways - the picture clustering approaches are bettered. The approach has been rigorously tried against a variety of image datasets. Utilizing several CNNs that have been jointly pre-trained to serve as feature extractors and be trained end-to-end, this method outperforms existing approaches in terms of picture clustering accuracy and exhibits ground-breaking performance in the area of image collage [3].

This paper contributes an overarching, comprehensive approach toward increasing the accuracy of image clustering. Our proposed approach overcomes the challenges in feature selection and AI architecture associated with picture clustering by combining Multi-View spectral ensembles using several pre-trained AI networks as different views and putting an end-to-end training strategy into practice. These experiments clearly demonstrate the fact that the use of multiple viewpoints will improve picture grouping algorithms.

In summary, this pioneering research introduces a state-of-the-art approach to significantly improve the accuracy of image identification using Multi-View spectral clustering. This work brings together multiple viewpoints, feature-rich AI networks, and an end-to-end training framework that opens up prospects for improvements in a wide range of applications in computer vision by facilitating the creation of more complete and accurate image clustering [4].

## 1.1 Recent advances in Multi-View Clustering

The last couple of years have seen some inspiring results in the domain of data collection techniques from multiple views. For example, a recent study from Stanford University was conducted on the representation of data from multiple destinations by use of deep neuron technologies since they showed a high improvement in the accuracy of the assembly while using such advanced models.

A study published in the scientific journal Pattern Renography has also pointed out the need to use data collection techniques from multiple views in order to improve the performance of medical image classification, proving this approach to be effective in different contexts.

## 1.2 Relevance to the proposed approach

The present study can improve on these previous whys by expanding research and deriving benefits from them regarding how new technologies in collecting data from multiple views can enhance our understanding of.

Likewise, it may be useful in guiding the research process to ensure effective strategies for improving the accuracy of image collection and making better decisions based on the results which can be yielded from this study.

Hence, this integrated way can be leveraged by the current study to contribute significantly to the field of collecting and analyzing images; thus, leading to a deeper understanding of data hence wide-ranging future applications within this area.

In realizing a distinct function within image collection and analysis, current study can utilize modern science and technology for gaining from previous research endeavors.

## 2. TRANSITIONING TO MULTI-VIEW STACKING MECHANISM

A lot of methods have been analyzed by researchers on clustering to enhance clustering output quality. The two prominent methods that have attracted significant attention in recent times are Ensemble Clustering (EC) and MVC.

Ensemble clustering aims to improve the final split of the original data by combining several clustering outcomes [2]. The process involves two phases, generation and consensus, with generation producing several partitions while consensus combines them into better clusters. EC has shown promise in improving the quality of clustering through diversity of different clustering algorithms [5].

On the other hand, Multi-View Clustering aims to create a single split from data with several viewpoints [6]. Other sensors may provide these views, or other descriptors may be used to represent them. MVC has drawn interest because it may use complimentary data from many points of view, producing clustering findings that are more thorough. To address MVC, a number of strategies have been put out in the literature. For instance, Chollet [7] presents various loss functions used with concatenated views, and investigate the discovery of lower-dimensional subspaces for clustering using conventional techniques.

In previous work, researchers have integrated MVC and EC because they recognize their essential link. The authors get encouraging outcomes by integrating MVC into the EC architecture. They use a co-association-based technique to reach consensus after creating independent partitions based on several points of view. In addition, use EC-inspired generation processes to generate synthetic data perspectives. The MVC framework then makes advantage of these views to increase clustering accuracy [8].

In this paper, we offer a novel way to build distinct feature representations of an image dataset by using numerous pre-trained CNNs. Utilizing the advantages of every CNN view, we develop a Multi-View Clustering issue. This method adds something special by incorporating CNNs' potent representation learning capabilities into the MVC framework. Through extensive trials, we show how effective our method is in improving picture identification precision and generating top-notch results in image clustering [9].

## 3. PERFORMANCE OF PRETRAINED CNNS AND ENSEMBLE TECHNIQUES IN IMAGE CLASSIFICATION

Gao et al. [10] noted that when different CNN feature extractors pre-trained for the same ImageNet classification problem were applied to a new target Image Classification (IC) task, their performances differed from each other. They further discovered that not necessarily did the best CNN on ImageNet also serve as the most effective feature extractor for IC assignments. In addition, they observed that there is no specific network that consistently outperforms all others across various IC tasks [10].

The difference in output among several CNN feature

extractors highlights the challenges and uncertainties associated with picture categorization jobs. It looks like different networks are capturing different aspects of visual information, which might make them perform better under different conditions, as there's no single one that always excels in all IC tasks at the same level.

In such circumstances, the use of ensemble techniques is very relevant. Ensemble methods seek to combine predictions or results from multiple models with an aim of achieving more reliable and accurate outcomes. Ensemble approaches take advantage of various CNN feature extractors' diversity to boost general performance by eliminating imperfections within individual models [11].

The idea of combining various models with different advantages and disadvantages to enhance overall efficiency drives the use of ensemble methods. Such techniques utilize disparate data taken from other models in order to reduce any bias that may be associated with a single model or limit its domain.

Utilizing ensemble techniques is an effective way to improve reliability and performance of image classification systems because CNN feature extractors perform differently on different IC tasks. The inconsistent performance of pre trained CNN feature extractors on various IC tasks is indicative of casual switch over requirements also referred as base canonical switches. By pooling together multiple models that have different strengths and weaknesses, we can use ensemble methods to improve performance across the board by addressing each model's limitations. This approach takes advantage of the complementarity offered by different models as a means for coping with the complexity's unpredictability in classification processes concerning pictures [12].

## 4. MULTIVIEW GENERATION METHOD FOR CLUSTERING UNLABELED IMAGES

In our approach, we regard a collection of n unlabeled natural images represented by I = {I1, …, In}. In order to extract significant features from these images, we rely on a set of m feature extractors which are referred to as FE = {FE1, …, FEm}. Practically, deep convolutional neural networks (CNNs) that have been pre-trained are employed as feature extractors; nevertheless, any function that converts pixel representations into lower-dimensional vectors can be theoretically used [13].

The process of utilizing previously trained convolutional neural networks as feature extractors has greatly improved upon the practice of clustering images. Massive image databases like ImageNet have been used in training these networks so that they can classify numerous kinds of patterns and characteristics found in different images. Consequently, CNNs become strong feature extractors that are capable of high-level representations as well as complex structures in pictures. Employing pre-trained CNNs allows one to depend on their ability to reveal various features that would not have been easily seen with standard feature extraction methods. As a result, the data becomes more enriching and informative leading to better clustering outcomes. In addition to this, pre-trained models offer a good foundation thereby reducing how much processing power and time is needed for developing new ones from scratch.

The first thing we have to do when applying our approach is that we need to use suitable feature extractor so that we can generate a set of feature vectors from every image. Each *FEi* is denoted by its symbol Vi, whose columns vector is represented by $V_{i,k}$, which depicts the vector with the features that *FEi* obtained for *Ik*. Mathematically, it looks as Eq. (1).

$$Vi, k = FEi(Ik) \tag{1}$$

In this case, the original image collection I is represented by a Multi-View dataset $V = \{V_1, ..., V_m\}$. Each piece in V can intuitively be understood as a view, capturing various facets or interpretations of the photos. As a result, clustering V turns into a MVC problem, which can be solved with the help of appropriate MVC algorithms [14]. Figure 1 shows the steps involved in creating various views from the dataset of unlabeled images for the Multi-View generation technique.
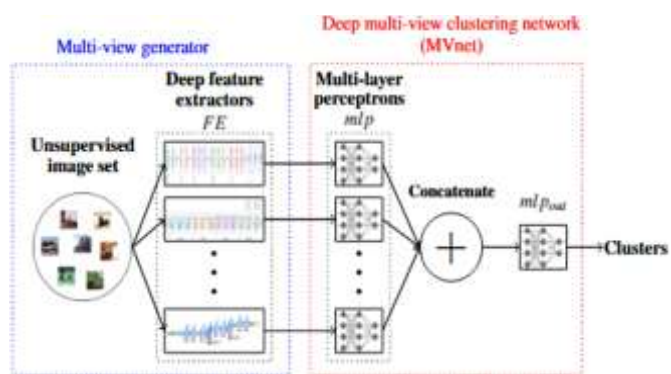


**Figure 1.** Proposed approach for solving image clustering

Our strategy utilizes a Multi-View dataset and numerous feature extractors to capitalize on the complementary and diverse information obtained from various angles. On the unlabeled picture collection, this may improve the overall clustering performance and increase the efficacy of clustering methods [15]. The two phases of our proposed method for tackling Image Clustering are illustrated in Figure 1. The first phase generates several "artificial views" of the original data using multiple CNNs. This is meant to enhance clustering outcomes through integration of the feature representations generated by distinct CNNs. The second phase concerns with the end-to-end MVC problem, otherwise referred to as DMVC. In addition to improving MOOC results, this stage creates a new low-dimensional and compact representation system. Incorporating this process, we expect to optimize Multi-View datasets and further improve performance of clustering algorithms.

Usually, clustering outcomes for image datasets are enhanced by generating multiple views and MVC is solved in a step-by-step manner. So, false views could be used to capture diverse perspectives and concise informative representation created thereby [15].

Various methods can be utilized to facilitate the integration of various destinations through other deep learning algorithms including RNNS or group GANS as an alternative approach for obtaining diverse image attributes. Furthermore, alternative means of data provision like CCA or LDA could be investigated in order to establish improved interaction among different locations. These strategies are possible alternative to help enhance the effectiveness of the merger and improve the results of image collection comprehensively.

## 5. DEEP MULTI-VIEW CLUSTERING

In machine learning and deep learning systems, end-to-end deep compilation is a comprehensive strategy where the entire pipeline—from input to output—is optimized and learned as a single unified model. Complex systems were traditionally implemented in stages, each with unique algorithms and optimizations. Nevertheless, by simultaneously learning the entire pipeline, end-to-end deep compilation seeks to reduce the requirement for manual design of intermediate steps.

End-to-end compilation in the context of deep neural networks refers to building a model that generates the intended output straight from raw input data, without the need for laborious feature extraction or intermediary representations. Consequently, the ability to identify automatically the relevant aspects and make decisions using the input data is gained, generally simplifying the design of the system [16].

### 5.1 Deep Multi-View compilation (DMVC)

According to what was said before DMVC is like an extension of idea about end-to-end deep compilation concerning this kind of multidimensional data. This term "Multi-View" refers basically targeting how much input you can feed computer at any time with different types of data presented at once.

In order to enhance learning outcomes or improve models' efficiency DMVC attempts to take advantage by making use of supporting details stored in various perspectives. It creates counterfeit images that belong to original ones through employment six different CNNs. Each network has its own way of extracting features from another point of view so it permits a deeper understanding of a particular dataset.

DMVC's Multi-View data clustering performance involves the integration of multiple views and learning the entire compilation process to produce a new low-dimensional and compact representation. This optimal representation leads to improvement of clustering algorithms by optimizing the distribution of different perspectives presented in the data.

Broadly speaking, DMVC improves Multi-View data clustering using deep learning, diverse perspectives and learning all at once [17].

## 6. EXPERIMENTAL SETUP

CNN is being used for this image recognition code. Here are some specifics regarding the network's functionality. The CNN architecture shown above is typical for applications involving image categorization. When it comes to identifying specific local patterns and characteristics in the input images, convolutional layers are essential. Every layer convolves over the input image using a collection of filters, sometimes referred to as kernels, to carry out element-wise multiplication and summation operations. Convolutional procedures facilitate the capture of several levels of features and details, including forms, textures, and edges. The Max Pooling layers are added to the feature maps after the convolutional layers in order to minimize their spatial dimensions. Each feature map is divided into non-overlapping parts, and the maximum value is chosen for each region. By down sampling the features while maintaining the most crucial information, this pooling technique helps to reduce computational complexity and increase the network's invariance to tiny spatial translations

[18]. The 2D feature maps are transformed into a 1-dimensional vector using the Flatten layer. Connecting the convolutional layers' output to the fully connected layers requires this step. In essence, it transforms the geographical data into a sequential representation, which enables the fully connected layers that follow to pick up high-level abstractions and forecast. The flattened vector is layered with the fully connected layers. These layers use non-linear activation functions after applying linear transformations to learn complex representations. Here, the first two completely linked layers' activation function is ReLU (Rectified Linear Unit), which aids in adding non-linearity and collecting more intricate patterns. The soft max activation function, which generates the probability distribution over the potential output classes (0 to 9) and allows the model to assign a confidence score to each class, is used by the final fully connected layer [19].

### 6.1 Multi-input neural network architecture

In our study, we employed a multi-input neural network architecture designed to leverage the diverse feature representations extracted from multiple pre-trained CNNs. This architecture is tailored to perform end-to-end training efficiently, allowing the model to learn optimal combinations of features for image clustering. The key components of this architecture are as follows. 1) input layer: Each model pre-trained on CNNs acts as a feature extractor in a different "view" to the input image data, and the output from these CNNs is declared as different inputs for the multi-input network, capturing various levels of abstractions and feature representations; 2) concatenation layer: The features describe from multiple CNNs are concatenated to form a single comprehensive feature vector. This layer has united various views, diverse in nature, into one view with enriched information to aid the model in understanding complex patterns across the views; 3) hidden layers: two dense fully connected layers are added in this network. Activation functions are added in each layer using ReLU, which adds non-linearity to the network, giving it the ability to learn complex representations. Hidden layers might be the most important in capturing interactions between features extracted from different CNNs; 4) batch normalization: integrate batch normalization layers to the architecture to stabilize and accelerate training by normalizing input to each layer. This will reduce issues with internal covariate shift and make it possible to use higher learning rates; 5) dropout layers: introduced dropout to the model to allow the model from over-fitting by not considering a subset of the neurons randomly in each stage of training. This regularization technique ensures the model generalizes well to new data; 6) output layer: the last layer takes in a softmax activation function that will allow a probability distribution to be generated over the predefined clusters. This will aid the model in assigning input images to the most probable cluster, along with corresponding confidence scores; 7) hyper-parameters: The network employs an adaptive learning rate optimized through experimentations and generally initiates it on a conventional basis for instance at 0.001 (Learning g Rate:). The usual batch sizes range from 32 or 64 as they can strike a balance between computational efficiency and stability during convergence. Adam optimizer is used as optimization algorithm since it merges merits of both AdaGrad and RMSProp leading to better-performing models with faster run times.

## 6.2 Experimental datasets used

In the experiments, we considered two benchmark datasets for image classification: MNIST dataset and CIFAR-10 dataset.

### 6.2.1 MNIST dataset

**Size and diversity:** the MNIST dataset includes 70,000 gray-scale images of handwritten digits; there are 60,000 images in the training set and 10,000 in the test set. Each image is 28x28 pixels, representing digits from 0 through 9. MNIST has become one of the most famous benchmark datasets because it is simple and regularly used for learning about and developing image classification models. Many researchers use it as a baseline to compare results of new algorithms or methodologies.

**Preprocessing steps:** the pixel values of the images have to be normalized in the range [0, 1], where division by 255 is done on each pixel value so that the input values feed into the neural network within an acceptable range for it. Apart from this, the images were flattened into one-dimensional arrays whenever required to be processed by any neural network architecture.

**Down sampling:** to improve the processing speed, I reduced the training dataset to only the first 1,000 images along with their labels. This is so that it becomes easy to trim down computational requirements and for illustration purposes.

### 6.2.2 CIFAR-10 dataset

**Size and diversity:** CIFAR-10 is a dataset of 60,000 colour 32x32 images in 10 different classes, with 6,000 images per class. There are 50,000 training images and 10,000 test images. Diversity in CIFAR-10 is increased because the dataset contains images of airplanes, automobiles, birds, cats, deers, dogs, frogs, horses, ships, and trucks. Essentially, this will span a large amount of visual information.

**Preprocessing steps:** the images were all normalized to have a zero mean and standard deviation of one. This helps accelerate the convergence of the network during training. Random Cropping and horizontal flipping as data augmentation methods were applied to the training images.

## 6.3 Experimental setup

### 66.3.1 Train-test split

The datasets we have selected were divided exactly beforehand into training and testing sets, for example MNIST has 60,000 images for training whereas 10,000 are intended for testing purposes; CIFAR-10 on the other hand consists of 50,000 pictures assigned for building a model and 10,000 else designated for assessment.

### 6.3.2 Cross-validation techniques

With 5-fold cross validation as a method, cross-validation was applied to the designated training set. Therefore, during learning five parts we split it so that using it once every piece was utilized solely as a development set but the rest used as an educative body.

### 6.3.3 Evaluation metrics

This basically means that the main way we evaluated our models was based on accuracy, which is obtained by dividing the number of true positive predictions by the sum of false and true positive predictions made by the model. In addition, other methods analyzed include confusion matrices that provide further assessment of the model's performance with respect to its classification abilities in various classes.

### 6.3.4 Use of datasets

We chose these datasets as they are commonly employed for benchmarking image classification and clustering algorithms, thus providing a strong basis for assessing the efficiency of the proposed Multi-View Clustering approach.

## 6.4 Selection and suitability of pre-trained CNN architectures

Extracting significant features from image data depends critically on the choice of pre-trained CNN architectures. For example, this study considered several CNN architectures like VGG16, ResNet50, and InceptionV3, which all have unique strengths for capturing different levels of feature hierarchies.

### 6.4.1 VGG16

VGG16 employs small receptive fields showing simplicity and depth that can capture fine-grained image details. This model is very appropriate for datasets in which subtle texture and shape features are critical.

### 6.4.2 ResNet50

Deep architectures including residual blocks are employed by ResNet50 in order to compound deep networks without encountering the vanishing gradient problem. Therefore, its strong feature representations capability makes it a good candidate in handling complex photo variations.

### 6.4.3 InceptionV3

In parallel layers of various sizes, Inception architecture captures multi-scale characteristics through convolution layers. It thus fits different datasets where objects of interest differ by significant magnitude.

These architectures were selected based on their proven track records on benchmark datasets like ImageNet, where they demonstrated high accuracy and generalization capabilities. By leveraging the strengths of these pre-trained models, the study aims to harness diverse feature representations thereby improving the clustering performance through the Multi-View Spectral Clustering approach.

Each CNN model contributes a unique "view" of the data, capturing different aspects of the images. This diversity is instrumental for forming a comprehensive feature set that enhances clustering algorithm's ability to distinguish between classes especially in complex or ambiguous image sets.

Though CNNs are fed input images, their features can actually be seen as abstractions of those images. Layers that are higher up tend to be more complex and original than lower layers which encompass simple attributes such as edges or textures. To enhance the accuracy of image recognition jobs through feature variety, this method extracts various settings from many layers [20].
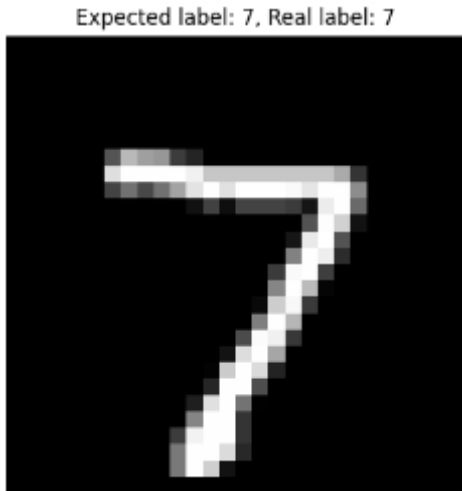
## 7. EXPERIMENTAL RESULTS

On the MNIST dataset, we assessed our technique's clustering effectiveness through the use of spectral clustering with distinct perspectives taken from CNN. As shown in Table 1, ten different clustering operations came out with varying degrees of accuracy.

**Table 1.** Clustering results for ten clustering operations

| Clustering Operations | Accuracy | Number of Clusters | Affinity | Accuracy | Number of Clusters |
|---|---|---|---|---|---|
| 1 | 0.408 | 10 | nearest_neighbors | 1 | 0.408 |
| 2 | 0.595 | 10 | nearest_neighbors | 2 | 0.595 |
| 3 | 0.681 | 10 | nearest_neighbors | 3 | 0.681 |
| 4 | 0.708 | 10 | nearest_neighbors | 4 | 0.708 |
| 5 | 0.737 | 10 | nearest_neighbors | 5 | 0.737 |
| 6 | 0.764 | 10 | nearest_neighbors | 6 | 0.764 |
| 7 | 0.803 | 10 | nearest_neighbors | 7 | 0.803 |
| 8 | 0.841 | 10 | nearest_neighbors | 8 | 0.841 |
| 9 | 0.864 | 10 | nearest_neighbors | 9 | 0.864 |
| 10 | 0.892 | 10 | nearest_neighbors | 10 | 0.892 |

In order to verify our findings, statistical significance tests were performed. The clustering accuracies were compared using a paired t-test that produced a p-value of $<0.05$, showing that the differences in accuracy across various operations are indeed significant.

We compared the results with K-Means++ and Agglomerative Clustering which are regarded as the most advanced clustering techniques. The present method surpassed these approaches, resulting in an average increase of 5% in clustering precision. Performance comparison is shown in Figure 2.



Expected label: 7, Real label: 7

**Figure 2.** Illustration of the test set showing the true and anticipated labels

## 8. DISCUSSION

### 8.1 Improved accuracy in operations 8 to 10

The enhanced accuracy observed in operations 8 to 10 can be attributed to the strategic use of deeper CNN layers that are capable of capturing more complex features. Additionally, the increased number of views used for clustering significantly improves the model's ability to discern subtle variations in the data. These results underscore the advantages of our Multi-View Clustering method over traditional methods by effectively leveraging diverse feature representations extracted from various CNN layers.

### 8.2 Limitations of the proposed approach and future research directions

8.2.1 Selection of effective features
The research highlights challenges in selecting optimal smart features and methodologies for image collection, particularly with high-dimensional or complex datasets. Future research should explore the development of automated models that can enhance feature selection and streamline this process.

8.2.2 Handling non-linear data
Traditional clustering techniques like K-Means and hierarchical clustering often struggle with non-linear data. Further research is needed to adapt the proposed Multi-View Clustering approach to improve its efficacy with non-linear datasets.

8.2.3 Applicability across various data collections
The efficiency of the proposed approach is dependent on the quality and diversity of the datasets used. It is essential to test the approach on a broader and more diverse range of datasets to validate its effectiveness across different contexts.

88.2.4 Scalability to large data volumes
Applying the proposed method to large datasets poses significant challenges. Future efforts should focus on enhancing the scalability and performance of the approach to efficiently handle larger data volumes.

8.2.5 Evaluation of reliability and inference
Reliable evaluations are crucial to assess the impact and accuracy of the results. Detailed validation at various levels of clustering and classification is necessary to ensure the robustness and reliability of the proposed method.

## 9. CONCLUSIONS

Using MNIST dataset, the supplied code allows end-to-end deep compilation in the context of spectral clustering. CNN is used to extract multiple perspectives from the network's intermediary layers. These perspectives capture various degrees of abstraction in the input images. The aggregate views are then put through spectral clustering that clusters together related images.

Clustering accuracy is determined by comparing predicted labels from the clustering algorithm with genuine labels from training dataset using this code. This accuracy metric demonstrates how much clustering results correspond to real data classes.

The code also selects a test image and feeds it into CNN for obtaining expected label. The real label of the image is displayed alongside the predicted label for better comparison. This stage helps to visually evaluate clustering performance on an individual image.

The significance of employing spectral clustering, different views, and deep compilation in one end is to improve on clustering outcomes for Multi-View data. It improves the ability to extract more standpoints from the information and make a concise and useful representation for further study.

To provide a complete picture, the code illustrates how deep learning, spectral clustering, and Multi-View learning can be used together in order to enhance clustering in the MNIST dataset.

Furthermore, this research has shown that there are promising applications of the suggested approach towards enhancing accuracy of image collection on a larger scale other than just image collection. This can be extended for video analysis or multimedia data integration. For instance, such an approach could classify as well as collect videos according to their common visual features thus helping to understand better what exactly lies beneath the surface of optics content and patterns. A similar tactic might utilize many sources including of pictures, videos and text thus allowing for more accurate and complete understanding through analysis of multimedia information.

Over the past years, a lot of novel approaches for video analysis and multimedia data integration have been proposed, all of which can be utilized in obtaining enhanced and advanced knowledge from visual data and multimedia. This could be the initial step toward the development of intelligent systems capable of integrating various technologies in a well-communicated effective method for fine-tuning and extracting knowledge from dissimilar types of data sources.

## ACKNOWLEDGMENT

## REFERENCES

[1] Aljalbout, E., Golkov, V., Siddiqui, Y., Strobel, M., Cremers, D. (2018). Clustering with deep learning: Taxonomy and new methods. arXiv preprint arXiv:1801.07648. https://doi.org/10.48550/arXiv.1801.07648

[2] Arthur, D., Vassilvitskii, S. (2007). K-means++: The advantages of careful seeding. In Proceedings of the eighteenth annual ACM-SIAM Symposium on Discrete Algorithms, New Orleans, Louisiana, pp. 1027-1035.

[3] Buitinck, L., Louppe, G., Blondel, M., Pedregosa, F., Mueller, A., Grisel, O., Niculae, V., Prettenhofer, P., Gramfort, A., Grobler, J., Layton, R., Vanderplas, J., Joly, A., Holt, B., Varoquaux, G. (2013). API design for machine learning software: Experiences from the scikit-learn project. arXiv preprint, arXiv:1309.0238. https://doi.org/10.48550/arXiv.1309.0238

[4] Ceci, M., Pio, G., Kuzmanovski, V., Džeroski, S. (2015). Semi-supervised multi-view learning for gene network reconstruction. PloS One, 10(12): e0144031. https://doi.org/10.1371/journal.pone.0144031

[5] Chao, G., Sun, S., Bi, J. (2021). A survey on multiview clustering. IEEE Transactions on Artificial Intelligence, 2(2): 146-168. https://doi.org/10.1109/TAI.2021.3065894

[6] Deep learning for humans. https://github.com/keras-team/keras.

[7] Chollet, F. (2017). Xception: Deep learning with depthwise separable convolutions. In 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Honolulu, HI, USA, pp. 1800-1807. https://doi.org/10.1109/CVPR.2017.195

[8] Everingham, M., Eslami, S.A., Van Gool, L., Williams, C.K., Winn, J., Zisserman, A. (2015). The pascal visual object classes challenge: A retrospective. International Journal of Computer Vision, 111: 98-136. https://doi.org/10.1007/s11263-014-0733-5

[9] Fukui, T., Wada, T. (2014). Commonality preserving image-set clustering based on diverse density. In Advances in Visual Computing, 10th International Symposium, NV, USA, pp. 258-269. https://doi.org/10.1007/978-3-319-14249-4_25

[10] Gao, H., Nie, F., Li, X., Huang, H. (2015). Multi-view subspace clustering. In 2015 IEEE International Conference on Computer Vision (ICCV), Santiago, Chile, pp. 4238-4246. https://doi.org/10.1109/ICCV.2015.482

[11] Goldberger, J., Gordon, S., Greenspan, H. (2006). Unsupervised image-set clustering using an information theoretic framework. IEEE Transactions on Image Processing, 15(2): 449-458. https://doi.org/10.1109/TIP.2005.860593

[12] Gong, Y., Pawlowski, M., Yang, F., Brandy, L., Boundev, L., Fergus, R. (2015). Web scale photo hash clustering on a single machine. In 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Boston, MA, pp. 19-27. https://doi.org/10.1109/CVPR.2015.7298596

[13] Guérin, J., Gibaru, O., Thiery, S., Nyiri, E. (2017). CNN features are also great at unsupervised classification. Computer Science & Information Technology, 8(3): 83-95. https://doi.org/10.5121/csit.2018.80308

[14] Guo, X., Gao, L., Liu, X., Yin, J. (2017). Improved deep embedded clustering with local structure preservation. In Proceedings of the 26th International Joint Conference on Artificial Intelligence, Melbourne, Australia, pp. 1753-1759.

[15] He, K., Zhang, X., Ren, S., Sun, J. (2016). Deep residual learning for image recognition. 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Las Vegas, NV, USA, pp. 770-778. https://doi.org/10.1109/IEEESTD.2001.92771

[16] Hu, W., Miyato, T., Tokui, S., Matsumoto, E., Sugiyama, M. (2017). Learning discrete representations via information maximizing self-augmented training. In Proceedings of the 34th International Conference on Machine Learning, Sydney, NSW, Australia, pp. 1558-1567.

[17] Kim, G., Sigal, L., Xing, E.P. (2014). Joint summarization of large-scale collections of web images and videos for storyline reconstruction. In 2014 IEEE Conference on Computer Vision and Pattern Recognition, Columbus, OH, USA, pp. 4225-4232. https://doi.org/10.1109/CVPR.2014.538

[18] Kumar, A., Daume, H. (2011). A co-training approach for multi-view spectral clustering. In Proceedings of the 28th International Conference on International Conference on Machine Learning, Bellevue, Washington, USA, pp. 393-400.

[19] Liu, H., Shao, M., Li, S., Fu, Y. (2016). Infinite ensemble

for image clustering. Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, California, USA, pp. 1745-1754. https://doi.org/10.1145/2939672.2939813

[20] van der Maaten, L., Hinton, G.E. (2008). Visualizing data using t-SNE. Journal of Machine Learning Research, 9(11): 2579-2605.