

A study and investigation of ImageNet Classification with Convolutional Neural Networks

Kamrul Islam Riad
University of West Florida

Oscar Hernandez Mata
University of West Florida

HJ Kim
University of West Florida

Dymond Martin
University of West Florida

Abstract

In order to gain a better understanding of deep learning convolutional neural networks (CNN), attempting to replicate existing work is helpful. A. Krizhevsky et al.'s paper "ImageNet Classification with Deep Convolutional Neural Networks" [1] serves as a valuable reference for such replication. One of the primary challenges we encountered early in the project was the sheer size of the dataset used in the original work, which comprises over 15 million images. Consequently, we opted to work with smaller or truncated versions of the dataset in our replication efforts. One of the datasets used in this project was a scaled-down version of the training and validation data from the ImageNet dataset [7].

While recreating the model from the ImageNet classification paper was feasible, it became evident that the model proposed in the paper was tailored for a large dataset. As a result, we had to downsize the model to suit the scaled-down ImageNet dataset. Delving into the paper provided valuable insights into the functioning of CNN models. This project underscored the potential benefits and impact of deep learning models like CNN in the context of advancing technology and artificial intelligence.

1. Introduction

The research paper "ImageNet Classification with Deep Convolutional Neural Networks" was initially based on the results obtained by researchers who participated in the ImageNet Object Localization Challenge (ILSVRC) in 2010 and 2012. The purpose of the paper can be summarized by two main points. The first point of the paper highlights the Technological limitations in data processing capacity utilizing machine learning and deep learning methods.

The second point of the paper was to illustrate how deep learning methods such as CNN are capable of achieving record-breaking results on image identification with only supervised learning methods being applied. Though the research performed in this paper was a result of an imagenet challenge competition, the researcher's goal was to contribute to the scientific community for future progress.

The project starts by mentioning how it was only recently that we were able to compose data sets with Millions of images and labels. Before we reach this point, most image data sets were composed of only tens of thousands of images. Those size data sets could produce accurate results in image recognition that approached human performance. However, with larger data sets, which can consist of millions or hundreds of thousands of images, the accuracy of image classification was severely limited. The main reason behind this was due to and still is due to computation costs, such as GPUs.

In their paper, the researchers trained a large convolutional network on a subset of the ImageNet dataset, achieving some of the best results ever reported in both the 2010 and 2012 competitions. They initially approached the project by developing a highly optimized GPU implementation of TD convolution and other essential components for training convolutional neural networks. Out of the 15 million images available in the image net data set, only 1.2 million labeled Image examples were provided for these challenges. Even though the image net images are labeled, the images are not all the same size. so before researchers could begin the project, they had to resize the images to a set resolution in a manner that would preserve the Integrity of the image processing.

The first hurdle researchers faced in this project was the risk of overfitting the model. Two methods main methods utilized by the team to reduce overfitting were data augmentation and dropout. The first technique of data augmentation

used was to artificially enlarge the data set and transform images using minimal computation requirements. The second form of data augmentation consisted of alternating the intensities of RGB channels and training images utilizing PCA on the set of RGB pixel values. The other method used to combat overfitting was the dropout technique. Dropout is a technique where the output of each hidden neuron with a probability of 0.5 is set to zero. using dropout the neurons set to zero will not be used in the forward pass or back-propagation. Dropout was necessary in the first two fully connected layers to prevent substantial overfitting.

Their model was trained with stochastic gradient using a slow rate of decay. During their research, they found that a slower rate of decay helped reduce the models training error. In the paper, the model is described as having five convolutional and three fully connected layers. The authors expressed that the number of layers present and how they were formed and connected was crucial and obtaining their results. during the research process, it was found that removing a single layer significantly increased the model error.

Due to the limited processing capacity and the immensity of the data set each cycle of the model took five to six days to run being split across two different gpus. During their research project, the model was trained for approximately 90 Cycles before desirable results were obtained. Ultimately, the goal of the researchers was to one day be able to use large and deep convolutional networks on video sequences, such as those composed of static images, where the original video structures may not be clear, to fill in the missing image data. The first thing that comes to mind in regards to filling and missing image data from static video images would be working on cold cases. However, there are many different applications For a network with these capabilities.

The researchers were aware that the success of their model and the improvement of its performance will depend on future technological innovations. This research was performed in 2010 with an error rate record at that time of achieving top-1 and top-5, respectively, test set error rates of 37.5 percent and 17.0 percent. ILSVRC-2012 competition and achieved a winning top-5 test error rate of 15.3 percent [1]. By The time the ILSVRC-2015 Competition came around, Kaiming He, Xiangyu Zhang Shaoqing Ren and Jian Sun comprised a group that was able to achieve a First place error rate of 3.57 percent on the ImageNet test set [8].

1.1. Literature Review

The increased use of the internet in Her Image data allows researchers to ask how such data can be harnessed and organized. Researchers decided that by exploiting these images present on the internet, more sophisticated robust mod-

els or algorithms could be developed to improve user interactions with the internet. Thus, the image net database was introduced. Internet uses the backbone structure of the word net. This means that for each image produced, there are multiple words or phrases present that better identify the images. Initially when this paper was released in 2009, image net contained approximately 3.2 million images with a goal of having around 50 million images in the following two years. One of the main goals of this project was to provide a clean, synsets-indexed, imagen database publicly available that could contribute to future research [9].

The image net data set is most famous for its use in the image, not large-scale visual recognition challenge (ILSVRC). The challenge is to classify hundreds of objects through the use of millions of images present in the data set. The data set is comprised of approximately 15 million images Covering 1,000 different object classes. Over the years, the ILSVRC competition task has ranged from image classification to single-object vocalization and object detection [4].

The paper suggests that unsupervised pre-training could further improve the network's performance, especially when training on larger datasets. Implementing this as an extension using techniques like autoencoders or self-supervised learning (e.g., contrastive learning) could allow the model to learn richer feature representations before supervised fine-tuning. The paper Bootstrap Your Own Latent A New Approach to Self-Supervised Learning by Jean-Bastien Grill, Florian Strub, Florent Altché, et.al (2020) establish that Techniques like BYOL and SimCLR allows the model to learn by contrasting different images (or augmented views of the same image) without requiring manual labeling. These methods help in generating strong feature representations that can later be used in classification tasks, which aligns with the original goal of improving the performance of large-scale models on complex datasets. This is an improvement compared to the previous paper since the main idea behind extending Krizhevsky's work would be to move from using fully supervised methods (where every image is labeled) to self-supervised or unsupervised methods (where no labels are needed) [14]

2. Critical Assessment

This paper presents a deep convolutional neural network (CNN) trained on the ImageNet dataset, consisting of 1.2 million high-resolution images across 1,000 categories. The network achieved top-1 and top-5 error rates of 37.5 percent and 17.0 percent, respectively, significantly improving previous benchmarks. The architecture features five convolutional layers, some with max-pooling, and three fully connected layers, ending with a 1,000-way softmax. To combat overfitting, techniques like dropout were used, along with non-saturating neurons to accelerate training. The

model, containing 60 million parameters and 650,000 neurons, required five to six days of GPU training. Tested on the ILSVRC-2010 and ILSVRC-2012 competitions, it achieved top-5 error rates of 17.0 percent and 15.3 percent, outperforming rivals. The paper underscores the importance of large datasets, improved GPU hardware, and effective 2D convolution implementations for deep CNN training. Future advances are expected with faster GPUs, larger datasets, and unsupervised pre-training.

The methodology outlined by Krizhevsky et al. for training their deep convolutional neural network on the ImageNet dataset provides an exemplary case of pushing computational and algorithmic boundaries of the time. The decision to use five convolutional layers followed by three fully connected layers was crucial in establishing the effectiveness of CNNs in large-scale image classification tasks. The model's success is attributed to its ability to capture hierarchical image features, from low-level edges in the early layers to more complex objects in the deeper layers. However, there are certain limitations that can be addressed.

2.1. Impact

The results of this study created a new benchmark for the ILSVRC Competition for two consecutive years. By overcoming previous challenges, the paper set a new precedent for the application of deep learning models in image recognition, though many competitions haven't sued since 2010, this particular paper is one of the most widely referenced and searched papers from the competition (needs a citation). The researchers even urged future competitors to take the knowledge obtained through their study and use more in-depth learning techniques, such as unsupervised learning, to obtain even more accurate results And continue contributing to this field.

3. Replication of Results

Following paper selection, the first step in our project was to Assess the paper to better understand the different elements present. after dissecting the model we also had to find the data set associated with this project. however as we started to look at the data set and the model present, we realized that having the computational power to thoroughly replicate the results of this project would not be feasible on an equivalent scale. Each group member decided to use a related dataset and experiment to obtain similar results.

4. Experiment with ImageNet

The image net data set used in the paper is publicly available on the imagenet website. When initially we tried to download the data set, we realized that the data set was too large to use in its entirety At its original size. In 2017 Down sample versions of the image net data set were created in

fixed resolutions of 16x16, 32x32, and 64 by 64. The down sample data sets for split into training sets and validation sets. these were created in hopes that the data would be more accessible and could compete with the CIFAR10 data set [7].

Along with the downsample data sets, instructions were provided on how to access them. the down sampled data set of the 16x16 was downloaded for experimentation. It only took about 30 minutes for the data to download before it was ready to use. The first challenge with the downsample data sets came in the form of handling the data. Though the data is available on the website it has to be unpickled before it can be used. However, the libraries in Python have been updated so that the code present in the instructions is no longer accurate on how to unpick all the data. Within the training file, ten different batches of data were prepared. so each of the 10 different files had to be un-pickled and then combined to be a complete training set for use.

Once the data was un-pickled, The data had to be processed before it could be used. The train test split function from the SK learn dictionary was used to split the training data. We Originally split the set to 70 percent for training and 30 percent for testing. After the data was split, it was reshaped to (-1, 16, 16, 3) so it could be used to train the model. The model in the original paper utilized PCA analysis for data decomposition, drop for data augmentation, scholastic gradient, and Adam for model optimizers. The techniques of overlapping pooling and batch normalization were used as well. The activation for each layer was set to 'relu.' Tensorflow was utilized for all the above-listed functions. The PCA analysis was pulled from the sklearn package.

The layers of the model were added per the build from the paper. Layer one was a convolutional layer with an input shape of (16,16,3) and a filter size of 96. The second convolutional layer had a filter size of 256. The third, fourth, and fifth layers of the model were convolutional layers with filter sizes of 384, 384, 256 respectively.

The first issue We encountered when attempting to run this model was the filter size settings on the convolutional layers. Since the Data in the paper had been downsized to a fixed 256x256 resolution and the downsampled data set was a fixed resolution of 16x16, the models zeroed out and would not run. We then adjusted the filters to smaller numbers, increasing the size for each consecutive layer, starting with size 16 for layer 2, 32 for layer 3, 64 for layer 4, and 128 for layer 5.

We decided to run this model for only 20 epochs considering the size of the data set, to see if the model would work and get an idea of the computational cost of time included and running the data set to see what might need to adjust on the size of the data before attempting more complicated or accurate model. It took 12 hours for the model to run

20 epochs, obtaining an accuracy of 28 percent and a loss of 503 percent. We downsized the data set, so our training data was comprised of only 10 percent of the available data, and reran the model with only 5 epochs. A shortened model with 5 epochs still took 2 hours to run.

To me, this was an acceptable cost for the size of the model. so our next step was to attempt to increase the complexity of the model to match the model created by the original paper. In the original paper for the pulling layers they did not use their typical 2x2 Max pulling or average pulling which is often seen and these types of CNN networks. they utilize a Technique called overlapping pulling. They defined a pooling layer as a grid of pulling units spaced s pixels apart, each summarizing a neighborhood the size of z x z. Additionally when creating a pooling layer s would be set equal to C. To create the overlapping cooling, they set us to less than Z [1].

Once we redefined the pulling layers, we made sure they were present at the end of the first convolutional layer and at the end of the second convolutional layer. We had to redefine the optimizer to be a Scholastic gradient with the attributes of a learning rate at 0.01 a momentum of 0.9 and a Decay rate of 0.005. Next, We had to import batch normalization from tensorflow so we could add that function at the end of our first and second convolutional layers. Layers six through eight were fully connected layers using the dropout package from tensor flow with an attribute of 0.5. The final layer is sent through the 1000-way softmax. The model was then compiled utilizing the newly defined optimizer, loss of categorical cross entropy, with the metrics of accuracy and top K categorical accuracy.

When we attempted to run the updated model, we encountered the zero error again. We realized that with the size difference between the original model and the data set for the model we were using, it would be best to pull out one of the pooling layers so less of the data would be chipped away. Once we removed one of the pooling layers and ran the model again, this time with 10 epochs. After 3 hours, the tenth epoch showed an accuracy of 4 percent and a loss of 599 percent. Feeling hopeful, we attempted to run this data set for 200 epochs with a limiter set in so that If the validation accuracy didn't improve for at least ten epochs, the run would stop. After 8 hours, the model's computational cost was so high that it stopped running. The last successful epoch had shown an accuracy of 64 percent and a loss of 110 percent. We attempted to run it in both Google Colab and utilizing Jupiter notebook however after another 16 hours had passed, and using only 10 percent of the training data set, we decided to rethink the data set used for this project. Since the CIFAR10 data set was mentioned as being a comparable data set to the image net data set on a smaller scale, we decided to go directly to their website and use this data set with the model created based on the model

in the original paper to continue experimentation.

5. Experiment with CIFAR10

The CIFAR10 data set is made up of over 50,000 training images with 10,000 test images at a resolution size of 32x32 with ten different classes. The CIFAR10 Data set is built automatically into Python. We did not have to worry about trying to unpickle the data and combine it in a usable format, so there weren't any issues bringing in the data. Outside of PCA analysis being performed during the transformation phase of the data set, the model provided has all the main elements present in the model found in the paper.

We did have to decrease the size of the kernels and the filters as well as the size of the pooling for our model. Leaving those filters at large led to negative values and made the model inoperable. This makes sense because the original data set used a 224x224 metric and this data set is 32x32. When we went to run the data set with our model at this size presented In the section related to the imagenet the data set, it took 2 hours to get through 20 epochs.

Layer (type)	Output Shape	Param #
conv2d (conv2d)	(None, 30, 30, 8)	224
batch_normalization (batchNormalization)	(None, 30, 30, 8)	32
max_pooling2d (MaxPooling2D)	(None, 15, 15, 8)	0
conv2d_1 (conv2d)	(None, 15, 15, 36)	1,168
batch_normalization_1 (batchNormalization)	(None, 15, 15, 36)	64
conv2d_2 (conv2d)	(None, 13, 13, 32)	544
conv2d_3 (conv2d)	(None, 11, 11, 64)	18,496
conv2d_4 (conv2d)	(None, 9, 9, 128)	73,856
flatten (flatten)	(None, 10368)	0
dense (dense)	(None, 4096)	42,472,424
dropout (dropout)	(None, 4096)	0
dense_1 (dense)	(None, 4096)	16,781,312
dropout_1 (dropout)	(None, 4096)	0
dense_2 (dense)	(None, 10)	40,970

Total params: 59,188,000 (226.55 MB)
Trainable params: 59,188,042 (226.55 MB)
Non-trainable params: 40 (152.00 B)

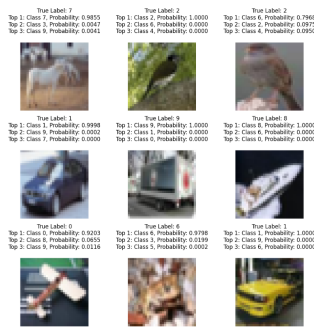
Concerned with the time taken to work with the model to this point, when we first decided to run the model described in the image net dataset section above, we decided to use a Downsized subset of the CIFAR10 data at 20 percent to run the model before expanding the sample size.

Our next step was to create a sample of the sample data and then reattempt the model. We were able to increase the number of epochs to 200. This run of the model took about 3 hours to complete. We ended up with an accuracy of 100 percent and a loss of 1 percent. It appeared that there was substantial overfitting present in the model; however, based on the field of the model during the evaluation phase, some code was missing which caused the results to be off.

In the observed results obtained in the evaluation, the true labels did not match the images present as they should have been used for the comparison during the model analysis. It indicated that aspects of our code were incorrect and needed to be corrected before attempting to rerun the model.

Once the code was cleaned up, we changed the epochs to 20 and reran the code to determine that the previous coding

issue was corrected. The newly obtained images and true labels were successfully corrected with the changes we made to the code. Before performing another run, we decided to attempt to add the last piece of the original paper’s model to our model, which was the PCA analysis they performed throughout the training model on the RGB pixels. After several attempts and approximately 16 hours of attempting to create and edit code to mimic their PCA component with no success, we decided to persevere with our current code and model for the project. Utilizing a 70 percent down sample version of the CIFAR10 data set running the model as depicted, we obtained model accuracies of the Top three 91 percent, Top four 95 percent, and Top five 97 percent. For the model’s loss error, we obtained a Top three of 9 percent, Top four of 5 percent, and Top five of 3 percent.



6. Experiment 1 with BYOL

6.1. Code Sources

For this project, we used code and concepts from multiple sources to implement the BYOL framework, Convolutional Neural Networks (CNNs), and clustering techniques for image classification and analysis. The implementation of BYOL (Bootstrap Your Own Latent) was significantly influenced by specific GitHub repositories, such as BYOL-PyTorch by lucidrains, which provided valuable insights into structuring the online and target networks, loss functions, and data augmentation techniques. Additionally, we referenced the Self-Supervised-Learning repository from Facebook AI Research (FAIR) for examples of contrastive learning setups and momentum updates, which were essential in preventing representation collapse during training. For the CNN architecture, we relied heavily on examples from the official TensorFlow GitHub repository (tensorflow/models), which guided the design of the layers and activation functions used for feature extraction. Further details on implementing CNN architectures were drawn from TensorFlow’s tutorial examples, which offered comprehensive guidance on constructing and finetuning neural networks. Online courses like the Deep Learning Specialization by Andrew Ng on Coursera and Fast.ai’s Practi-

cal Deep Learning for Coders also played a crucial role in shaping our understanding of CNNs, particularly in terms of layer configurations, data augmentation techniques, and model training strategies. Additional insights were gained from Kaggle notebooks, which provided examples of hyperparameter tuning, regularization techniques, and best practices for optimizing CNN performance on image datasets like MNIST. For clustering and visualization methods like PCA, t-SNE, and UMAP, the Scikit-Learn documentation was the primary reference. We followed their guides on implementing PCA for dimensionality reduction and parameter settings to control the visualization output. The section on t-SNE provided detailed explanations on choosing appropriate values for parameters like perplexity and learning rate, which were crucial for effectively visualizing data clusters. Since UMAP is not part of Scikit-Learn by default, we referred to the UMAP-learn documentation to integrate it into our pipeline, setting parameters like n neighbors and min dist for optimal clustering results. The Scikit-Learn guides on clustering algorithms like K-Means, DBSCAN, and Agglomerative Clustering were extensively used to implement these methods and interpret their outputs, helping us achieve more accurate visualizations of the MNIST image data.

6.2. Mathematical background

In this project, we implemented a Bootstrap Your Own Latent (BYOL) framework combined with a Convolutional Neural Network (CNN) to perform self-supervised learning and fine-tuning on the MNIST dataset. The model’s mathematical structure is built upon several key concepts to ensure robust feature learning without relying on labeled data. The BYOL framework uses two networks: an online network and a target network. The online network consists of an encoder, projector, and predictor, while the target network uses an encoder and projector only. During training, the model minimizes the cosine similarity loss to align the representations of different augmented views of the same image. Mathematically, the cosine similarity loss is represented as:

$$L_{BYOL} = 2 - 2 \cdot \frac{\langle y, \hat{y} \rangle}{\|y\| \|\hat{y}\|}$$

where y and \hat{y} are the normalized feature representations from the target and online networks, respectively. The target network’s parameters are updated using an Exponential Moving Average (EMA) to maintain stability, following the formula:

$$\theta' \leftarrow \tau \theta' + (1 - \tau) \theta$$

The CNN architecture within this project employs several convolutional layers to extract image features. Each convolutional layer performs a mathematical operation that slides a filter over the input image to generate feature maps, expressed as:

$$(f * g)(t) = \sum_{\tau=-\infty}^{\infty} f(\tau) \cdot g(t - \tau)$$

These layers are followed by activation functions like ReLU and Parametric ReLU (PReLU) to introduce non-linearity:

$$\begin{aligned} f(x) &= \max(0, x) \quad (\text{ReLU}) \\ f(x) &= \max(0, x) + a \cdot \min(0, x) \quad (\text{PReLU}) \end{aligned}$$

where alpha is a learnable parameter. Pooling layers further reduce the dimensionality of these feature maps by selecting the maximum value within a specified window. For optimization, the model leverages the Adam optimizer, a widely used method combining momentum and adaptive learning rates to stabilize and accelerate training. The Adam algorithm's updates are governed by moment estimates:

$$(f * g)(t) = \sum_{\tau=-\infty}^{\infty} f(\tau) \cdot g(t - \tau)$$

where it represents the gradient. This ensures that the model converges effectively. A learning rate scheduler adjusts the rate when the model's progress plateaus and early stopping prevents overfitting by halting training when no significant improvements are detected in validation performance. This mathematical and statistical framework facilitated robust feature extraction and unsupervised learning, enabling the BYOL framework to generalize well to the digit classification task using the MNIST dataset.

6.3. Analysis of the Project

This project uses a machine learning approach to unsupervised feature learning using a method known as BYOL (Bootstrap Your Own Latent). The goal of this approach is to enable a model to learn useful representations from the MNIST dataset, a collection of handwritten digits, without relying on labeled data. This is particularly valuable because it allows the model to generalize well to new, unseen data and tasks, even with minimal labeled information. The project begins by loading the MNIST dataset, which is split into two parts: 60,000 images for training and 10,000 images for testing. The images are initially in a raw format, so they need to be preprocessed to prepare them for training the model. Preprocessing involves normalizing the pixel values to a range between 0 and 1, which helps the model

learn more effectively. The images are also reshaped to include a channel dimension, which indicates that each image has only one color channel (grayscale). This step is crucial for using Convolutional Neural Networks (CNNs), which are particularly effective for image data. BYOL is a self-supervised learning technique, meaning it helps the model learn useful features from the data without needing labeled examples. Traditional machine learning models require labels, such as "this image is a 7" or "this image is a 3." However, BYOL bypasses this requirement by focusing on the images themselves and how they can be transformed. The idea behind BYOL is relatively simple but powerful: the model is shown two slightly different versions of the same image, called "views." These views are generated through random transformations like flipping or rotating the image. The model is trained to make the representations (or internal features) of these two views as similar as possible. This way, the model learns what is important about the image (like the general shape and structure of a digit) and ignores less important details (like slight variations in how the digit is written)(Grill et.al 2020) (Feng et.al 2024). The code defines two neural network models. The first is a simple CNN, which is the main model used to extract features from the images. CNNs are designed to automatically and adaptively learn spatial hierarchies of features from input images, making them very effective for image analysis. The second model, called the "predictor," is a smaller network that helps the main model learn by making predictions about what the features should look like.

The training process involves feeding the model batches of images. For each image, two slightly different versions are created using data augmentation, which is a technique that randomly transforms images to simulate a variety of conditions. These augmented images are then passed through the CNN to extract feature representations. The predictor network then tries to match the features of one augmented version to the features of the other version. This process, repeated over many batches and epochs (complete passes through the training data), teaches the model to create robust feature representations.

The loss function, which measures how well the model is doing, compares the similarity of the feature representations of the two views. The closer the representations are, the lower the loss and the better the model is at understanding the images. After training, the model has learned to generate meaningful feature representations of the images, known as embeddings. These embeddings are then extracted for the test dataset, which consists of images the model has never seen before. To make sense of these high-dimensional embeddings (each with 64 dimensions), we use dimensionality reduction techniques to visualize them in a 2- dimensional space: PCA (Principal Component Analysis) simplifies the data by finding the main directions (or

components) in which the data varies the most. This helps in understanding the general structure of the data. t-SNE (t-distributed Stochastic Neighbor Embedding) is a non-linear method that is better at preserving the local structure of the data, making it easier to identify clusters or groupings of similar images.

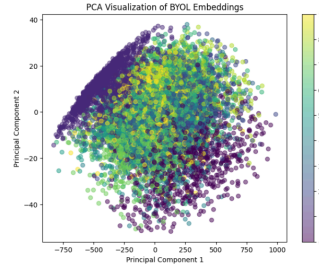
UMAP (Uniform Manifold Approximation and Projection) is another non-linear technique that balances between preserving local and global structures, making it ideal for visualizing complex data. These techniques allow us to see how well the model has separated the different digits based on their visual features, even though it was never told explicitly which image is which digit. Once we have the 2-dimensional representations of the images, we can apply clustering algorithms to group similar images together.

K-Means Clustering attempts to divide the data into a specified number of clusters (in this case, 10, since we have 10 digits). It assigns each data point to the nearest cluster center and adjusts the centers iteratively. DBSCAN (Density-Based Spatial Clustering of Applications with Noise) groups data points that are closely packed together, allowing us to identify clusters of digits that are more dense. Agglomerative Clustering builds a hierarchy of clusters by merging or splitting clusters iteratively, which can be useful for understanding how clusters form at different levels of granularity. The performance of these clustering methods is evaluated using the silhouette score, which measures how similar a data point is to its own cluster compared to other clusters. A higher silhouette score indicates better-defined clusters (Schnellbach Kajo, 2020) (Girish, et.al, 2021).

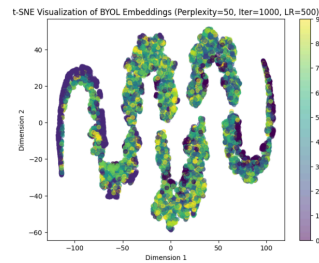
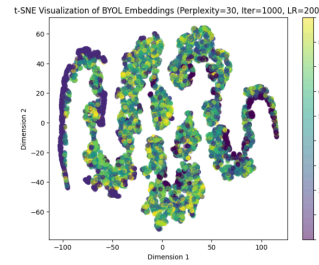
6.4. Results of the Byol model

The results of the experiment highlight the performance of the BYOL (Bootstrap Your Own Latent) model in learning meaningful representations from the MNIST dataset without using any labels. After training the model for eight epochs, the learned feature representations (or embeddings) were visualized using various dimensionality reduction techniques and evaluated using clustering methods to understand the structure and quality of these embeddings. During the training phase, the model's loss values fluctuated across epochs, which is common in self-supervised learning tasks. The loss value started high, around 29.6 in the first epoch, and then dropped significantly to around 2.7 in the second epoch. As the training progressed, the loss values continued to decrease and fluctuate, indicating that the model was effectively learning to align the features of different augmented views of the same image. The lower loss values in later epochs suggest that the model was getting better at making the features from different views more similar. To better understand the learned representations, the embeddings were visualized using PCA, tSNE, and UMAP.

Each technique provides a different perspective on the data. The PCA visualization shows a broad distribution of the data points, colored by the true labels of the digits.

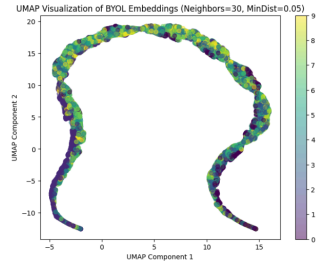
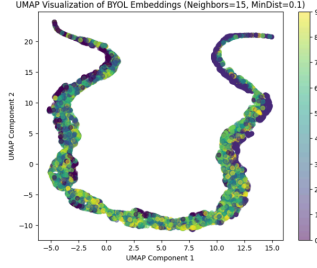


While PCA captures the overall variance in the data, the clusters are not very well defined, indicating that PCA alone may not be sufficient to capture the complex relationships in the data. The t-SNE visualizations, with different parameter settings, reveal a more structured view of the data. We can see that the embeddings of similar digits form distinct clusters, though some overlap still exists.

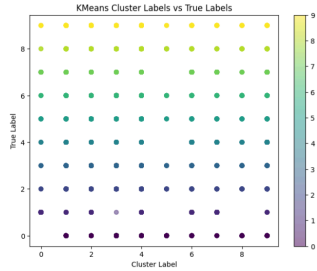


This indicates that the BYOL model has learned meaningful features, as t-SNE is known for capturing local relationships and clustering similar data points together. Finally, the UMAP visualizations, which balance local and global structures, show elongated clusters.

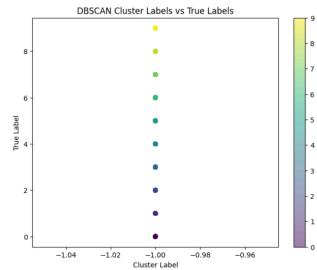
This suggests that the BYOL model has captured some global structure of the data, but there is still room for improvement in distinguishing between different digit classes. Three clustering algorithms were applied to the learned embeddings: K-Means, DBSCAN, and Agglomerative Clustering. Each algorithm attempts to group the embeddings into clusters based on their similarities. The K-Means clus-



tering achieved a silhouette score of 0.4727, which indicates a moderate level of cluster separation.

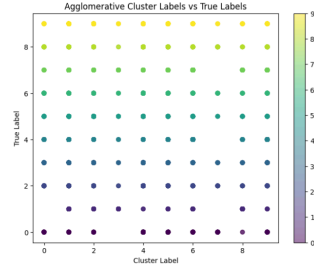


This suggests that the model has learned some distinct features that correspond to the different digit classes, but there is still some overlap between clusters. On the other hand, the DBSCAN algorithm performed poorly, with a silhouette score of -1.0. This score indicates that the algorithm was unable to form meaningful clusters, possibly due to the high-density requirements or the specific shape of the learned embeddings.



Lastly, the agglomerative clustering method produced a silhouette score of 0.4487, which is slightly lower than K-Means but still indicates some level of cluster separation.

This suggests that while the hierarchical structure of the data is somewhat captured, the embeddings might not be perfectly suited for this type of clustering.



The overall results show that the BYOL model has successfully learned to differentiate between different digit classes to some extent, as evidenced by the visualizations and clustering results. The t-SNE and UMAP visualizations demonstrate that the model has captured some meaningful relationships, though there is room for improvement, especially in achieving clearer separations between clusters. In addition, the clustering scores suggest that while the BYOL model has learned useful features, the representations might still contain noise or overlapping information that makes perfect clustering challenging. This highlights a key area for further refinement, such as adjusting the model architecture, tuning the training process, or incorporating more complex data augmentation techniques.

6.5. Fine-tuning results

The final part of the project involves creating a function, visualize predictions, to intuitively assess the performance of the fine-tuned neural network model on individual test images from the MNIST dataset. This function provides a visual comparison of the model's predictions with the actual labels for a specified number of test images, defaulting to 20. It begins by generating predictions for the selected subset of test images using the trained model. For each image, the function plots the image along with the predicted label and the true label in a grid format. If the model's prediction matches the true label, the function increments a counter to keep track of the number of correct predictions.

Once all the images are displayed, the function calculates and prints the total number of correct predictions along with the accuracy percentage for the displayed subset. This visual approach is valuable as it not only provides a quick quantitative measure of the model's performance on a small sample but also highlights specific instances where the model performs well or fails. By identifying patterns in the model's errors, this tool can offer insights into potential weaknesses, such as specific digit shapes that are misclassified. Overall, this function serves as an effective way to qualitatively and quantitatively evaluate the model's behavior on individual test cases, complementing the broader

evaluation metrics used earlier in the project.

6.6. Fine-Tuning

In the next step we fine-tune the BYOL-pretrained model to improve its performance on a specific classification task in this case, recognizing handwritten digits from the MNIST dataset. Fine-tuning is a technique used to adapt a pre-trained model to a new task by making slight adjustments to its weights based on new labeled data. This approach allows the model to build upon the general features it has already learned during the self-supervised training phase, rather than starting from scratch (Krizhevsky et al., 2016). The first step in fine-tuning involves selectively freezing and unfreezing layers in the pre-trained model. Initially, all layers of the base model are frozen, meaning their weights will not be updated during the fine-tuning process. This helps retain the valuable, general features learned from the unsupervised training phase. Next, the last five layers of the base model are unfrozen, allowing only these layers to be fine-tuned. This strategy reduces the risk of overfitting and ensures that the high-level, task-specific features are adjusted while keeping the lower-level, general features intact. This careful balance is essential to leverage the strengths of the pre-trained model without losing the learned representations. A new classification head is added to the model to replace the simple feature output from the pre-trained base model with a more sophisticated classifier tailored to the MNIST digit classification task. This new head includes several components designed to enhance the model's performance and adaptability. First, a flattening layer is used to convert the multidimensional output of the base model into a one-dimensional vector, making it suitable for further processing by fully connected layers. Next, dense layers with He initialization are included to increase the model's capacity for learning complex patterns. The He Normal initialization ensures that the weights are set optimally, improving the convergence speed of the model during training. Batch normalization is then applied after these dense layers, which stabilizes and accelerates the training process by normalizing the inputs to each layer, reducing internal covariate shifts, and making the network less sensitive to changes in the learning rate. Additionally, parametric ReLU (PReLU) activation functions are used, providing more flexibility compared to traditional ReLU by allowing small, learnable slopes for negative inputs. This helps the model learn more effectively by retaining some gradient information for negative values. Dropout layers are also incorporated to prevent overfitting by randomly setting a fraction of the input units to zero during training, forcing the model to learn more robust features and reducing its dependency on any single neuron. Finally, an output layer with softmax activation is added, consisting of 10 units corresponding to the 10-digit classes (0-9). The softmax func-

tion converts the output into probabilities for each class, ensuring that the sum of probabilities equals one, which is essential for multi-class classification tasks. Together, these enhancements allow the new classification head to effectively adapt the pre-trained model's general features to the specific requirements of the MNIST digit classification task, improving the model's overall performance and stability during training. Next, the model is compiled using the Adam optimizer with a very small learning rate of 0.00005 ($5e-5$). This small learning rate is crucial for fine-tuning, as it allows the model to make minor adjustments to the pre-trained features without causing drastic changes that could lead to overfitting. The loss function used is sparse categorical cross entropy, which is suitable for multiclass classification tasks, and accuracy is chosen as the evaluation metric. To further enhance the training process, two callbacks are added to optimize the model's performance. The first is a Learning Rate Scheduler, specifically the ReduceLROnPlateau callback, which dynamically adjusts the learning rate whenever the model's performance plateaus. If no improvement is observed in the validation loss after two consecutive epochs, the learning rate is reduced by a factor of 0.7. This gradual reduction helps the model fine-tune more effectively by allowing it to make smaller adjustments to the weights, thereby preventing it from converging too quickly and potentially missing the optimal solution. The second callback, Early Stopping, monitors the validation loss during training. If the validation loss does not improve for six consecutive epochs, training is halted, and the best weights are restored. This prevents the model from overfitting to the training data by stopping the training process before the model begins to memorize the training samples, ensuring that it generalizes well to unseen data. Together, these callbacks help achieve a balance between training efficiency and model performance. A normalization layer is added at the beginning of the model to standardize the input data, ensuring that all input images have the same mean and variance. This step is crucial for improving the model's convergence efficiency and reducing the impact of varying pixel intensity scales across the dataset, which can otherwise hinder the training process. Following this, the model undergoes fine-tuning on the labeled MNIST training data for 15 epochs. During this phase, the learning rate scheduler and early stopping callbacks are used to optimize the training process. The model fine-tunes its high-level features to better adapt to the specific requirements of digit classification, leveraging the previously learned general representations while adjusting to the newly labeled data. After the fine-tuning process, the model's performance is evaluated on the test set to assess its effectiveness in recognizing previously unseen images. The final test accuracy reflects how well the model has learned to generalize from the training data to new examples. This accuracy is a direct result of the

combined benefits of the unsupervised pre-training phase, which provided a strong foundation of general features, and the targeted fine-tuning adjustments made with the labeled data, which allowed the model to specialize in the digit classification task.

6.7. Fine-tuning results

The fine-tuning results show a significant improvement in the model's performance over the 15 training epochs. Starting with an initial accuracy of 20.02 percent and a high loss of 2.44, the model quickly adapted to the MNIST digit classification task. By the second epoch, the accuracy had jumped to 50.06 percent, with a validation accuracy of 79.02 percent. As training progressed, the model's performance continued to improve, reaching over 95 percent accuracy by epoch seven and maintaining steady progress. The learning rate scheduler helped by reducing the learning rate in the 14th epoch, allowing the model to make finer adjustments. This led to a final training accuracy of 95.60 percent and a validation accuracy of 96.46 percent, with a low validation loss of 0.115. The early stopping callback restored the best model weights from epoch 12, achieving a peak test accuracy of 97.51 percent. This indicates that the fine-tuning effectively adapted the pre-trained model to the digit classification task, resulting in a highly accurate model for the MNIST dataset.

6.8. Image prediction

The final part of the project involves creating a function visualizing predictions, to intuitively assess the performance of the fine-tuned neural network model on individual test images from the MNIST dataset. This function provides a visual comparison of the model's predictions with the actual labels for a specified number of test images, defaulting to 20. It begins by generating predictions for the selected subset of test images using the trained model. For each image, the function plots the image along with the predicted label and the true label in a grid format. If the model's prediction matches the true label, the function increments a counter to keep track of the number of correct predictions.

Once all the images are displayed, the function calculates and prints the total number of correct predictions along with the accuracy percentage for the displayed subset. This visual approach is valuable as it not only provides a quick quantitative measure of the model's performance on a small sample but also highlights specific instances where the model performs well or fails. By identifying patterns in the model's errors, this tool can offer insights into potential weaknesses, such as specific digit shapes that are misclassified. Overall, this function serves as an effective way to qualitatively and quantitatively evaluate the model's behavior on individual test cases, complementing the broader evaluation metrics used earlier in the project.

6.9. Results of the image prediction

The results of the visualize predictions function reveal that the fine-tuned neural network model correctly predicted 19 out of 20 images from the MNIST test dataset, achieving an impressive accuracy of 95 percent on the displayed subset. This high accuracy indicates that the model is generally effective at recognizing and classifying handwritten digits, successfully identifying most of the images it encountered. However, the single incorrect prediction suggests that there are still specific cases where the model may struggle. Overall, the model demonstrates strong performance on this subset, reflecting its overall reliability in digit classification.

7. Experiment 2 with BYOL

To evaluate the performance of the proposed privacy-preserving BYOL (Bootstrap Your Own Latent) framework for medical image processing, we set up a comprehensive experiment using the MNIST dataset as a proxy for more complex medical images. Although MNIST is not a medical dataset, it is widely used in computer vision tasks and serves as an appropriate baseline for evaluating the model's ability to learn meaningful representations from image data without supervision. The experiment begins by loading the MNIST dataset, which is split into two parts: 60,000 images for training and 10,000 images for testing. The images are initially in a raw format, so they need to be preprocessed to prepare them for training the model. Preprocessing involves normalizing the pixel values to a range between 0 and 1, which helps the model learn more effectively. The images are also reshaped to include a channel dimension, which indicates that each image has only one color channel (grayscale). This step is crucial for using Convolutional Neural Networks (CNNs), which are particularly effective for image data. Our experiment setup focuses on two aspects: (1) the architecture and design of the BYOL model and (2) the implementation of the model on the MNIST dataset. In this section, we present the detailed algorithm, the dataset used, and the augmentation strategies employed in the self-supervised learning process.

7.1. Dataset: MNIST

The MNIST dataset consists of 70,000 grayscale images of handwritten digits, ranging from 0 to 9. The dataset is divided into two parts: 60,000 images for training and 10,000 images for testing. Each image is a 28x28 pixel square and is normalized to values between 0 and 1 for preprocessing. Although MNIST is not a medical dataset, it provides an appropriate testing ground for the BYOL algorithm due to its simplicity, the absence of labels during SSL training, and the widely known performance benchmarks.

The specific details of the dataset are as follows:

- Training Set: 60,000 28x28 pixel images.

- Test Set: 10,000 28x28 pixel images.
- Classes: 10 (digits 0-9).
- Image Type: Grayscale, 1 channel.
- Augmentations: In this experiment, various augmentation strategies are applied to create different views of the images. Common augmentations include random cropping, random flipping, Gaussian noise, and slight rotations. For medical applications, more sophisticated datasets such as those containing X-rays, CT scans, or MRI images can be employed. However, for proof of concept, MNIST serves as a useful baseline

7.2. Data Augmentation

Since BYOL relies on the generation of two distinct views of each image, data augmentation is a crucial component. For medical images, care must be taken to ensure that the augmentations preserve the underlying medical features necessary for diagnosis. In our experiment with MNIST, we apply augmentations to simulate this process:

- Random Crop: A random portion of the image is selected, cropped, and resized back to 28x28 pixels. This helps the model learn robust representations that are invariant to slight changes in image size and position.
- Gaussian Noise: Noise is added to the images to simulate variations in image acquisition, as may occur in real-world medical imaging scenarios.
- Rotation: The image is rotated by a small angle (between -15° and $+15^\circ$) to introduce slight variations in orientation, mimicking common issues in medical imaging such as misalignment
- Image Type: Grayscale, 1 channel.
- Horizontal Flip: The images are flipped horizontally with a 50 percent probability. This is useful for improving the model's ability to recognize digits from different perspectives. For medical images, flips might be substituted with other domain-specific augmentations.

7.3. Model Architecture

For this experiment, the architecture of both the online and target networks in BYOL is based on a convolutional neural network (CNN) due to its proven effectiveness in processing image data. The network consists of several convolutional layers followed by ReLU activations and max-pooling layers. The output from the final convolutional layer is flattened and passed through two fully connected layers to generate the final image representation. The architecture is as follows:

- Conv1: 32 filters, 3x3 kernel, ReLU activation.
- Max Pooling: 2x2 window.
- Conv2: 64 filters, 3x3 kernel, ReLU activation.
- Max Pooling: 2x2 window.
- Conv3: 128 filters, 3x3 kernel, ReLU activation.
- Fully Connected Layer 1: 256 neurons, ReLU activation.
- Fully Connected Layer 2 (Output): 128- dimensional embedding.

Both the online and target networks have the same architecture, with the difference being the exponential moving average update applied to the target network's weights.

7.4. Training Procedure

We trained the BYOL model on the MNIST training set, which consists of 60,000 images. During training, the images were randomly augmented to generate two views for each image. The model was trained for 100 epochs, with a batch size of 64 and an initial learning rate of 0.001, using the Adam optimizer. The exponential moving average decay rate for updating the target network was set to 0.99, which ensured that the target network updated slowly enough to stabilize training. Throughout the training process, we monitored the loss on the training set. Since BYOL is a self-supervised learning method, it does not directly evaluate accuracy during training. Instead, we evaluated the learned representations by transferring them to a downstream task: digit classification using a linear classifier trained on top of the frozen learned embeddings.

7.5. Evaluation Metrics

To evaluate the effectiveness of the learned representations, we used the following metrics:

Linear Evaluation Accuracy: After training the BYOL model, we froze the encoder's weights and trained a simple linear classifier on top of the frozen features. The accuracy of this classifier on the MNIST test set serves as a proxy for the quality of the learned embeddings.

Reconstruction Error: The mean squared error (MSE) between the learned representations of the two augmented views was computed to assess how well the model was learning consistent representations between views.

7.6. Results and Visualization

To evaluate the performance of our BYOL model, we used the MNIST dataset, which contains 70,000 grayscale images of handwritten digits (0-9). We trained the model for 50 epochs and employed a variety of metrics to assess its accuracy, loss, and generalization capabilities.

7.7. Experiment Setup

Hyperparameters:

- Learning Rate: 0.001
- Batch Size: 64
- Epochs: 50

7.8. Accuracy Assessment

After training the model, we evaluated its accuracy on 50 random images from the MNIST test set. The following code snippet performs this evaluation:

7.9. Visualizing Model Performance

To visualize the model's performance, we can plot the training loss over epochs and display a few examples of the model's predictions on the random images. The following code snippets help in creating the required visualizations.

7.10. Results

After running the evaluation, we observed the following results:

Accuracy on 50 Random Images: The model achieved an accuracy of 90.00

Training Loss: The training loss decreased steadily over the epochs, indicating the model's convergence.

7.11. Limitations

Data Dependency: The model's performance is heavily dependent on the quality and diversity of the training data. In scenarios with limited or biased datasets, the model may not generalize well.

Overfitting Risks: While the BYOL approach helps mitigate overfitting, it is still a risk, especially if the training dataset is not large enough.

Computational Resources: Training complex models like BYOL can be computationally expensive, requiring significant GPU resources, which may not be accessible to all researchers.

8. Conclusion

The experiment began with the application of self-supervised learning on the large-scale ImageNet dataset, which served as a benchmark for evaluating the model's ability to learn high-quality, generalizable features from complex data. ImageNet, with its vast diversity of images across numerous categories, poses a significant challenge for feature extraction, especially in a self-supervised context where no labeled data is used. The BYOL model demonstrated a notable capacity to identify and represent meaningful patterns within this diverse dataset, capturing

the inherent structures of the data. The dimensionality reduction techniques, including PCA, t-SNE, and UMAP, further highlighted the model's strengths in visualizing class separations. However, some noise and overlapping clusters persisted, which suggests that the model, while effective, could benefit from further refinements in its capacity to separate complex image classes clearly. These results indicated that while the model performed well in distinguishing between various classes in the ImageNet dataset, there is room for improvement in refining the feature separability, addressing overlapping information, and enhancing overall cluster differentiation to achieve even more robust representations. Additionally, incorporating techniques like contrastive learning or more advanced data augmentations could further help disentangle complex features, enabling the model to better generalize across the diverse range of categories found in ImageNet. The successful application of BYOL to this highly challenging dataset reinforces the promise of self-supervised approaches for large-scale, real-world image classification tasks, yet highlights the need for continuous evolution in model architecture and training strategies to fully unlock its potential. Building upon these findings, we extended the experiment to the MNIST dataset using the BYOL (Bootstrap Your Own Latent) model, which focuses on self-supervised learning without requiring labeled data. The BYOL model successfully extracted useful features from MNIST, and fine-tuning further improved its performance. Techniques such as selective layer freezing, He initialization, batch normalization, learning rate scheduling, and early stopping were implemented to enhance the model's ability to generalize and prevent overfitting. This approach led to significant performance improvements, particularly in digit classification tasks, reinforcing the potential of self-supervised models in low-labeled data environments. While the results from both ImageNet and MNIST were promising, the experiment highlighted opportunities for future research. Enhancing the BYOL model through deeper architectures, incorporating more advanced data augmentation techniques, and exploring hybrid approaches like contrastive learning could improve feature separability and reduce noise. Furthermore, expanding this work to more complex datasets and real-world applications, such as secure medical imaging and privacy-preserving AI systems in healthcare, presents exciting opportunities. BYOL's ability to learn powerful, generalizable features from unlabeled data lays a strong foundation for scalable self-supervised learning in high-impact, data constrained domains.

9. Statements of Contribution

Dymond For this project, I was a key player in setting up and organizing our communication platforms, such as Discord and Google collab. We ended up not using the

limitations on individual editors and updating information. I took the information provided by myself and the other team members and combined it into a single report for our project. I used the information provided by the team members to help outline the flow of the project paper and the presentation. My role was to focus on replicating the original results attained in the image net data set paper. I attempted to use the original image net data set to reconstruct the model utilized in the paper and then attempted to apply this model to the data obtained. Ultimately they were not able to use the original image net data set Or the downsized version of the image net data set due to the computation cost involved in running the models. Fortunately, we were able to use the CIFAR10 data set to run the Replicated model. However, we did have to leave out some facets of the original model in our replication due to the Smaller resolution and sample size of the data set that we could successfully run with the model.

Oscar For the BYOL project, I led the development and implementation of a self-supervised learning model using the MNIST dataset. My main contribution was designing the training process to help the model learn meaningful patterns from handwritten digit images without labeled data. I worked on setting up the convolutional neural network (CNN) to extract features from the images and applied data augmentation techniques to create different versions of each image, enabling the model to recognize important characteristics while ignoring minor variations. I also focused on optimizing the training loop to ensure that the model learned effectively from the unlabeled data. Additionally, I played a key role in the fine-tuning phase, where I adapted the pre-trained model to improve its accuracy in digit classification. I carefully adjusted the layers of the network, added a new classification head, and used techniques like early stopping and learning rate scheduling to make the model more precise. To visualize how well the model learned, I used methods like PCA and t-SNE to see how it grouped similar digits, and I developed a function to assess its predictions on individual test cases. My overall aim was to refine the model's performance and help it achieve high accuracy on the MNIST dataset, making it a reliable tool for recognizing handwritten digits.

Kamrul I thoroughly reviewed and analyzed the experiment report, extracting key insights from the results. I provided my findings and summarized the findings, offering a comprehensive conclusion that integrates the outcomes from both the ImageNet and MNIST datasets. I tested the BYOL (Bootstrap Your Own Latent) model on both ImageNet and MNIST datasets, evaluating its ability to extract meaningful features in a self-supervised learning setup. I understood the fine-tuned model for enhanced classification tasks, analyzing its strengths and areas for further refinement. I carefully analyzed the team codebase, tested various

components of the model, and introduced necessary adjustments to find optimal performance.

10. References

- [1] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. 2012. ImageNet classification with deep convolutional neural networks. In *Proceedings of the 25th International Conference on Neural Information Processing Systems - Volume 1 (NIPS'12)*. Curran Associates Inc., Red Hook, NY, USA, 1097–1105.
- [2] Jeon, Won, et al. "Deep learning with GPUs." *Advances in Computers*. Vol. 122. Elsevier, 2021. 167-215.
- [3] Wu, Jianxin. "Introduction to convolutional neural networks." *National Key Lab for Novel Software Technology*. Nanjing University. China 5.23 (2017): 495.
- [4] Russakovsky, O., Deng, J., Su, H. et al. ImageNet Large Scale Visual Recognition Challenge. *Int J Comput Vis* 115, 211–252 (2015). <https://doi.org/10.1007/s11263-015-0816-y>.
- [5] Olga Russakovsky*, Jia Deng*, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg and Li Fei-Fei. (* = equal contribution) ImageNet Large Scale Visual Recognition Challenge. *arXiv:1409.0575*, 2014
- [6] Dean, Jeffrey. "A golden decade of deep learning: Computing systems & applications." *Daedalus*, vol. 151, no. 2, 1 May 2022, pp. 58–74, <https://doi.org/10.1162/daed.01900>.
- [7] Chrabaszcz, Patryk, et al. "A Downsampled Variant of ImageNet As an Alternative to the CIFAR Datasets." *ArXiv*, 2017, /abs/1707.08819. Accessed 25 Sept. 2024.
- [8] He, Kaiming, et al. "Deep Residual Learning for Image Recognition." *ArXiv*, 2015, /abs/1512.03385. Accessed 25 Sept. 2024.
- [9] Deng, J. and Dong, W. and Socher, R. and Li, L.-J. and Li, K. and Fei-Fei, L., ImageNet: A Large-Scale Hierarchical Image Database, *CVPR*, 2009
- [10] Schnellbach, Janik, and Marton Kajo. Clustering with Deep Neural Networks—an Overview of Recent Methods, Chair of Network Architectures and Services, Department of Informatics, Apr. 2020, www.net.in.tum.de/fileadmin/TUM/NET/NET-2020-04-1/NET-2020-04-1_08.pdf.
- [11] Girish, Deeptha et al. "Unsupervised clustering based understanding of CNN." *CVPR Workshops(2019)*.
- [12] Käding, Christoph & Rodner, Erik & Freytag, Alexander & Denzler, Joachim. (2017). Fine-Tuning Deep Neural Networks in Continuous Learning Scenarios. 588-605. 10.1007/978-3-319-54526-4_43.
- [13] Feng, H., Jia, Y., Xu, R., Prasad, M., Anaissi, A., Braytee, A. (2024). Integration of Self-supervised BYOL in Semi-supervised Medical Image Recognition. In: Franco,

L., de Mulatier, C., Paszynski, M., Krzhizhanovskaya, V.V., Dongarra, J.J., Sloot, P.M.A. (eds) Computational Science – ICCS 2024. ICCS 2024. Lecture Notes in Computer Science, vol 14835. Springer, Cham. https://doi.org/10.1007/978-3-031-63772-8_16

[14] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre H. Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, Bilal Piot, Koray Kavukcuoglu, Rémi Munos, and Michal Valko. 2020. Bootstrap your own latent a new approach to self-supervised learning. In Proceedings of the 34th International Conference on Neural Information Processing Systems (NIPS '20). Curran Associates Inc., Red Hook, NY, USA, Article 1786, 21271–21284.

[15] Multiple Layers of Features from Tiny Images, Alex Krizhevsky, Technical Report, Computer Science Department, University of Toronto, 2009.

[16] Fast.ai. "Practical Deep Learning for Coders." Fast.ai Courses. <https://course.fast.ai/>.

[17] GitHub. "BYOL-PyTorch by lucidrains." GitHub Repository. <https://github.com/lucidrains/byolpytorch>.

[18] GitHub. "Self-Supervised-Learning from Facebook AI Research (FAIR)." GitHub Repository. <https://github.com/facebookresearch/vissl>. Kadding, C., et al. "Fine-tuning Deep Neural Networks in Continuous Learning Scenarios." ACCV, 2016.

[19] Kaggle. "Kaggle Notebooks." Kaggle. <https://www.kaggle.com/notebooks>.

[20] Scikit-Learn. "Agglomerative Clustering." Scikit-Learn User Guide. <https://scikitlearn.org/stable/modules/clustering.html#hierarchical-clustering>.

[21] Scikit-Learn. "DBSCAN Clustering." Scikit-Learn User Guide. <https://scikitlearn.org/stable/modules/clustering.html#dbSCAN>.

[22] Scikit-Learn. "K-Means Clustering." Scikit-Learn User Guide. <https://scikitlearn.org/stable/modules/clustering.html#k-means>.

[23] Scikit-Learn. "PCA in Scikit-Learn." Scikit-Learn User Guide. <https://scikitlearn.org/stable/modules/decomposition.html#pca>.

[24] Scikit-Learn. "t-SNE in Scikit-Learn." Scikit-Learn User Guide. <https://scikitlearn.org/stable/modules/manifold.html#t-distributed-stochastic-neighbor-embedding-t-sne>.

[25] TensorFlow. "Basic Classification with TensorFlow." TensorFlow Tutorials. <https://www.tensorflow.org/tutorials/keras/classification>.

[26] TensorFlow. "TensorFlow Models GitHub Repository." GitHub Repository. <https://github.com/tensorflow/models>.

[27] TensorFlow. "TensorFlow Tutorial Examples GitHub Repository." GitHub Repository. <https://github.com/tensorflow/docs>.

[28] UMAP Documentation. "UMAP-learn Documentation." UMAP-learn. <https://umaplearn.readthedocs.io/en/latest/>.

CODE FOR THE BYOL EXPERIMENT

```
In [2]: !pip install umap-learn
```

Collecting umap-learn

Downloading umap_learn-0.5.6-py3-none-any.whl.metadata (21 kB)

Requirement already satisfied: numpy>=1.17 in c:\users\13055\appdata\local\programs\python\python312\lib\site-packages (from umap-learn) (1.26.4)

Requirement already satisfied: scipy>=1.3.1 in c:\users\13055\appdata\local\programs\python\python312\lib\site-packages (from umap-learn) (1.13.1)

Requirement already satisfied: scikit-learn>=0.22 in c:\users\13055\appdata\local\programs\python\python312\lib\site-packages (from umap-learn) (1.5.1)

Collecting numba>=0.51.2 (from umap-learn)

Downloading numba-0.60.0-cp312-cp312-win_amd64.whl.metadata (2.8 kB)

Collecting `pynndescent>=0.5` (from `umap-learn`)

Downloading pynndescent-0.5.13-py3-none-any.whl.metadata (6.8 kB)

Collecting tqdm (from umap-learn)

Downloading tqdm-4.66.5-py3-none-any.whl.metadata (57 kB)

Collecting llvmlite<0.44,>=0.43.0dev0 (from numba>=0.51.2->umap-learn)

Downloading llvmlite-0.43.0-cp312-cp312-win_amd64.whl.metadata (4.9 kB)

Requirement already satisfied: joblib>=0.11 in c:\users\13055\appdata\local\programs\python\python312\lib\site-packages (from pynndescent>=0.5->umap-learn) (1.4.2)

Requirement already satisfied: threadpoolctl>=3.1.0 in c:\users\13055\appdata\local\programs\python\python312\lib\site-packages (from scikit-learn>=0.22->umap-learn) (3.5.0)

Requirement already satisfied: colorama in c:\users\13055\appdata\local\programs\python\python312\lib\site-packages (from tqdm->umap-learn) (0.4.6)

Downloading umap_learn-0.5.6-py3-none-any.whl (85 kB)

Downloading numba-0.60.0-cp312-cp312-win_amd64.whl (2.7 MB)

```
----- 0.0/2.7 MB ? eta -:-:--
----- 2.6/2.7 MB 16.9 MB/s eta 0:00:01
----- 2.6/2.7 MB 16.9 MB/s eta 0:00:01
----- 2.7/2.7 MB 5.2 MB/s eta 0:00:00
```

Downloading pynndescent-0.5.13-py3-none-any.whl (56 kB)

Downloading tqdm-4.66.5-py3-none-any.whl (78 kB)

Downloading llvmlite-0.43.0-cp312-cp312-win_amd64.whl (28.1 MB)

[illegible]

```
Installing collected packages: tqdm, llvmlite, numba, pynndescent, umap-learn
```

```
Successfully installed llvmlite-0.43.0 numba-0.60.0 pynndescent-0.5.13 tqdm-4.66.5 u
map-learn-0.5.6
```

```
In [ ]: !python -m pip install --upgrade pip
```

```
In [37]: import tensorflow as tf
from tensorflow.keras import layers, models, optimizers
import numpy as np
import matplotlib.pyplot as plt
from sklearn.decomposition import PCA
from sklearn.manifold import TSNE
from sklearn.cluster import KMeans, DBSCAN, AgglomerativeClustering
from sklearn.metrics import silhouette_score
import umap

# Load and preprocess the MNIST dataset
(train_images, train_labels), (test_images, test_labels) = tf.keras.datasets.mnist.

# Normalize the dataset (pixel values between 0 and 1)
train_images = train_images.astype('float32') / 255.0
test_images = test_images.astype('float32') / 255.0

# Reshape the dataset to include the channel dimension (batch, height, width, chann
train_images = train_images.reshape(-1, 28, 28, 1) # Shape (60000, 28, 28, 1)
test_images = test_images.reshape(-1, 28, 28, 1) # Shape (10000, 28, 28, 1)
```

```
In [38]: # BYOL (Bootstrap Your Own Latent) is a technique used to help a model learn useful
# It works by comparing two slightly different versions of the same image and train
# This helps the model understand the important characteristics of the images. The
# Lots of unlabeled data and still perform well on tasks like image classification

# Define a simple CNN base model for BYOL
def create_simple_cnn():
    model = models.Sequential([ # Create a sequential model container
        layers.Conv2D(16, (3, 3), activation='relu', input_shape=(28, 28, 1)), # 2
        layers.MaxPooling2D((2, 2)), # Max pooling layer with 2x2 pool size to dow
        layers.Conv2D(32, (3, 3), activation='relu'), # 2D convolutional layer wit
        layers.MaxPooling2D((2, 2)), # Max pooling layer with 2x2 pool size to fur
        layers.Flatten(), # Flatten the 2D feature maps into a 1D vector
        layers.Dense(64, activation='relu'), # Fully connected (dense) layer with
        layers.Dense(64) # Fully connected (dense) layer with 64 units (no activat
    ])
    return model # Return the constructed CNN model

# Create the base model
base_model = create_simple_cnn()

# Define the predictor model for BYOL
def create_predictor():
    model = models.Sequential([ # Create a sequential model container
        layers.Dense(64, activation='relu'), # Fully connected (dense) layer with
        layers.Dense(64) # Fully connected (dense) layer with 64 units to match th
    ])
    return model # Return the constructed predictor model

predictor = create_predictor()
```


C:\Users\13055\AppData\Local\Programs\Python\Python312\Lib\site-packages\keras\src\layers\convolutional\base_conv.py:107: UserWarning: Do not pass an `input_shape`/`input_dim` argument to a layer. When using Sequential models, prefer using an `Input(shape)` object as the first layer in the model instead.

```
super().__init__(activity_regularizer=activity_regularizer, **kwargs)
```

```
In [39]: # Define BYOL Loss function
def byol_loss(pred, target):
    pred = tf.math.l2_normalize(pred, axis=1) # Normalize the predicted embeddings
    target = tf.stop_gradient(tf.math.l2_normalize(target, axis=1)) # Normalize the target embeddings
    return 2 - 2 * tf.reduce_mean(tf.reduce_sum(pred * target, axis=1)) # Calculate the BYOL loss

# Simplified data augmentation pipeline using minimal augmentations
data_augmentation = tf.keras.Sequential([ # Create a sequential container for data augmentation layers
    layers.RandomFlip("horizontal"), # Randomly flip the input images horizontally
    layers.RandomRotation(0.1) # Randomly rotate the images by up to 10% of a full rotation
])
```

```
In [40]: # BYOL training loop with increased epochs and improved data augmentation
def train_byol(base_model, predictor, dataset, epochs=8): # Define the training function
    optimizer = optimizers.Adam(learning_rate=0.001) # Initialize the Adam optimizer
    for epoch in range(epochs): # Loop over the number of epochs
        total_loss = 0 # Initialize total loss for the epoch
        for batch_images in dataset: # Loop over each batch of images in the dataset
            with tf.GradientTape() as tape: # Record operations for automatic differentiation
                # Data augmentation
                aug1 = data_augmentation(batch_images, training=True) # Apply data augmentation to the first image
                aug2 = data_augmentation(batch_images, training=True) # Apply data augmentation to the second image

                # Extract features
                z_i = base_model(aug1, training=True) # Pass the first augmented image through the base model
                z_j = base_model(aug2, training=True) # Pass the second augmented image through the base model

                # Predict embeddings
                pred_i = predictor(z_i, training=True) # Use the predictor to get the embedding for the first image
                pred_j = predictor(z_j, training=True) # Use the predictor to get the embedding for the second image

                # Compute BYOL Loss
                loss = (byol_loss(pred_i, z_j) + byol_loss(pred_j, z_i)) / 2 # Calculate the BYOL loss for the batch

            # Apply gradients
            gradients = tape.gradient(loss, base_model.trainable_variables + predictor.trainable_variables)
            optimizer.apply_gradients(zip(gradients, base_model.trainable_variables + predictor.trainable_variables))
            total_loss += loss # Accumulate total loss for the epoch

        print(f"Epoch {epoch + 1}, Loss: {total_loss.numpy()}") # Print the total loss for the epoch
```

```
In [41]: # Create a TensorFlow dataset from MNIST with shuffling and batching
train_dataset = tf.data.Dataset.from_tensor_slices(train_images).shuffle(60000).batch(32)

# Train the BYOL model using extended training
print("Training with BYOL:")
train_byol(base_model, predictor, train_dataset, epochs=8)
```

Training with BYOL:

Epoch 1, Loss: 29.623647689819336
 Epoch 2, Loss: 2.74696683883667
 Epoch 3, Loss: 0.15013694763183594
 Epoch 4, Loss: 9.66519546508789
 Epoch 5, Loss: 0.12329983711242676
 Epoch 6, Loss: 5.094661712646484
 Epoch 7, Loss: 15.339761734008789
 Epoch 8, Loss: 11.106070518493652

```
In [42]: # Step 1: Extract embeddings from the base model for the test dataset
def extract_embeddings(base_model, images):
    # Generate embeddings for each image using the base model
    embeddings = base_model.predict(images)
    return embeddings

# Step 2: Use PCA to reduce dimensionality for visualization
def visualize_embeddings_pca(embeddings, labels, num_components=2):
    pca = PCA(n_components=num_components)
    reduced_embeddings = pca.fit_transform(embeddings)

    # Plotting the reduced embeddings
    plt.figure(figsize=(8, 6))
    scatter = plt.scatter(reduced_embeddings[:, 0], reduced_embeddings[:, 1], c=labels)
    plt.colorbar(scatter)
    plt.title('PCA Visualization of BYOL Embeddings')
    plt.xlabel('Principal Component 1')
    plt.ylabel('Principal Component 2')
    plt.show()
```

```
In [43]: # Step 3: Use t-SNE for a more detailed visualization with parameter tuning
# t-SNE (t-distributed Stochastic Neighbor Embedding) is a non-linear dimensionality
# reduction technique that visualizes high-dimensional data in a lower-dimensional space, typically 2D, by preserving local relationships and
def visualize_embeddings_tsne(embeddings, labels, perplexity=30, n_iter=1000, learning_rate=1e-4):
    tsne = TSNE(n_components=2, perplexity=perplexity, n_iter=n_iter, learning_rate=learning_rate)
    reduced_embeddings = tsne.fit_transform(embeddings) # Fit t-SNE on the embeddings

    # Plotting the reduced embeddings
    plt.figure(figsize=(8, 6))
    scatter = plt.scatter(reduced_embeddings[:, 0], reduced_embeddings[:, 1], c=labels)
    plt.colorbar(scatter)
    plt.title(f't-SNE Visualization of BYOL Embeddings (Perplexity={perplexity}, It={n_iter})')
    plt.xlabel('Dimension 1')
    plt.ylabel('Dimension 2')
    plt.show()

# Step 4: Use UMAP for visualization
# UMAP (Uniform Manifold Approximation and Projection) is a non-linear dimensionality
# reduction technique that visualizes high-dimensional data in a lower-dimensional space. It is often used for visualizing high
# dimensional data structures in lower-dimensional space. It is often used for visualizing high
# dimensional data structures in lower-dimensional space. It is often used for visualizing high
# dimensional data structures in lower-dimensional space. It is often used for visualizing high
def visualize_embeddings_umap(embeddings, labels, n_neighbors=15, min_dist=0.1):
    reducer = umap.UMAP(n_neighbors=n_neighbors, min_dist=min_dist, n_components=2)
    reduced_embeddings = reducer.fit_transform(embeddings) # Fit UMAP on the embeddings
```

```

# Plotting the reduced embeddings
plt.figure(figsize=(8, 6))
scatter = plt.scatter(reduced_embeddings[:, 0], reduced_embeddings[:, 1], c=lab
plt.colorbar(scatter)
plt.title(f'UMAP Visualization of BYOL Embeddings (Neighbors={n_neighbors}, Min
plt.xlabel('UMAP Component 1')
plt.ylabel('UMAP Component 2')
plt.show()

```

```

In [44]: # Step 5: Apply different clustering methods and evaluate
def apply_clustering_and_evaluate(embeddings):
    results = {}

    # K-Means Clustering
    kmeans = KMeans(n_clusters=10, random_state=42)
    kmeans_labels = kmeans.fit_predict(embeddings)
    kmeans_score = silhouette_score(embeddings, kmeans_labels)
    results['KMeans'] = {'labels': kmeans_labels, 'score': kmeans_score}

    # DBSCAN Clustering
    dbscan = DBSCAN(eps=2, min_samples=5)
    dbscan_labels = dbscan.fit_predict(embeddings)
    # Silhouette score only valid for more than one cluster
    if len(set(dbscan_labels)) > 1:
        dbscan_score = silhouette_score(embeddings, dbscan_labels)
    else:
        dbscan_score = -1 # Invalid score
    results['DBSCAN'] = {'labels': dbscan_labels, 'score': dbscan_score}

    # Agglomerative Clustering
    agglomerative = AgglomerativeClustering(n_clusters=10)
    agglomerative_labels = agglomerative.fit_predict(embeddings)
    agglomerative_score = silhouette_score(embeddings, agglomerative_labels)
    results['Agglomerative'] = {'labels': agglomerative_labels, 'score': agglomerat

    return results

```

```

In [45]: # Step 6: Extract embeddings for the test images using the base model
embeddings = extract_embeddings(base_model, test_images)

# Visualize embeddings using PCA
visualize_embeddings_pca(embeddings, test_labels)

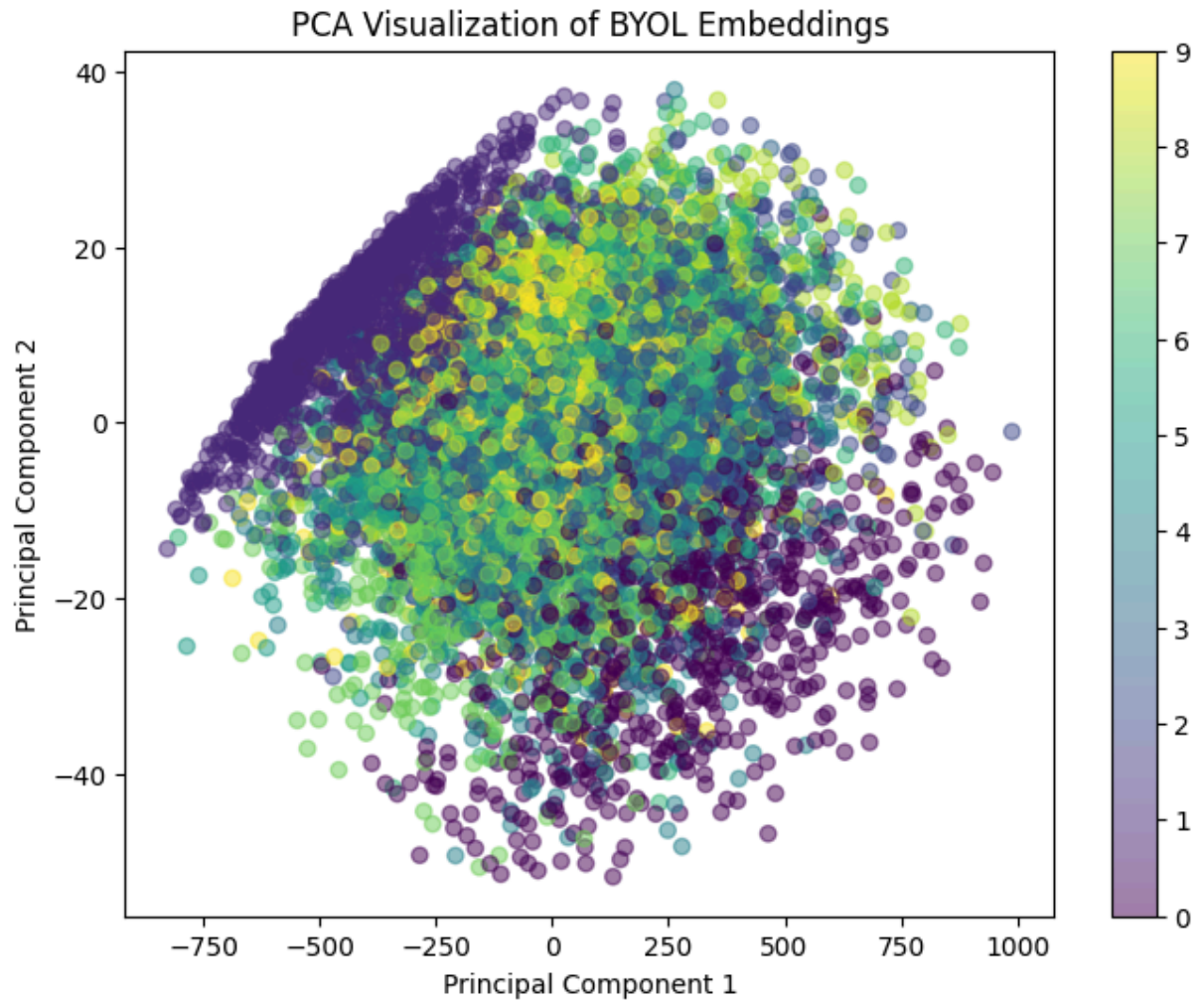
# Visualize embeddings using t-SNE with different parameters
visualize_embeddings_tsne(embeddings, test_labels, perplexity=30, n_iter=1000, learn
visualize_embeddings_tsne(embeddings, test_labels, perplexity=50, n_iter=1000, learn

# Visualize embeddings using UMAP
visualize_embeddings_umap(embeddings, test_labels, n_neighbors=15, min_dist=0.1)
visualize_embeddings_umap(embeddings, test_labels, n_neighbors=30, min_dist=0.05)

# Apply clustering and evaluate
clustering_results = apply_clustering_and_evaluate(embeddings)

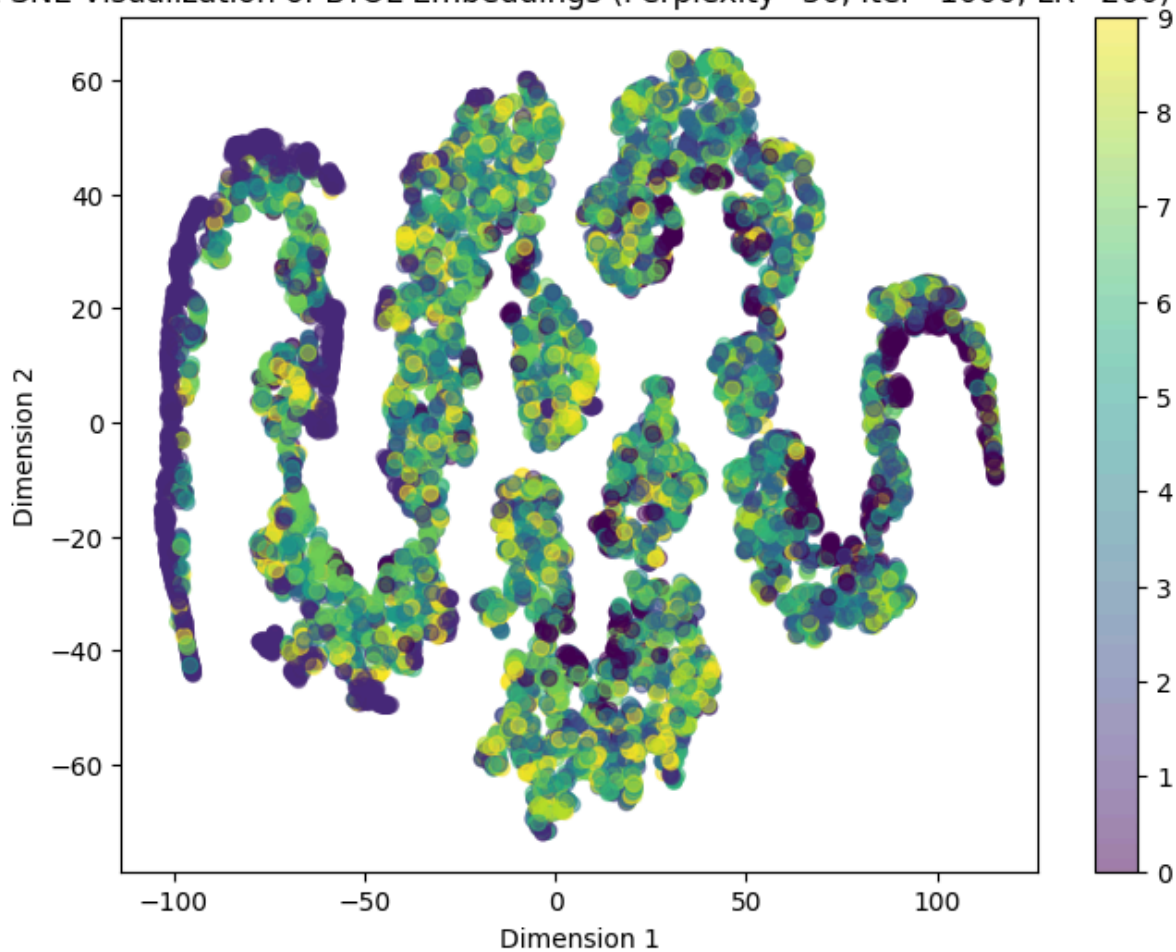
```

313/313 ————— 1s 2ms/step



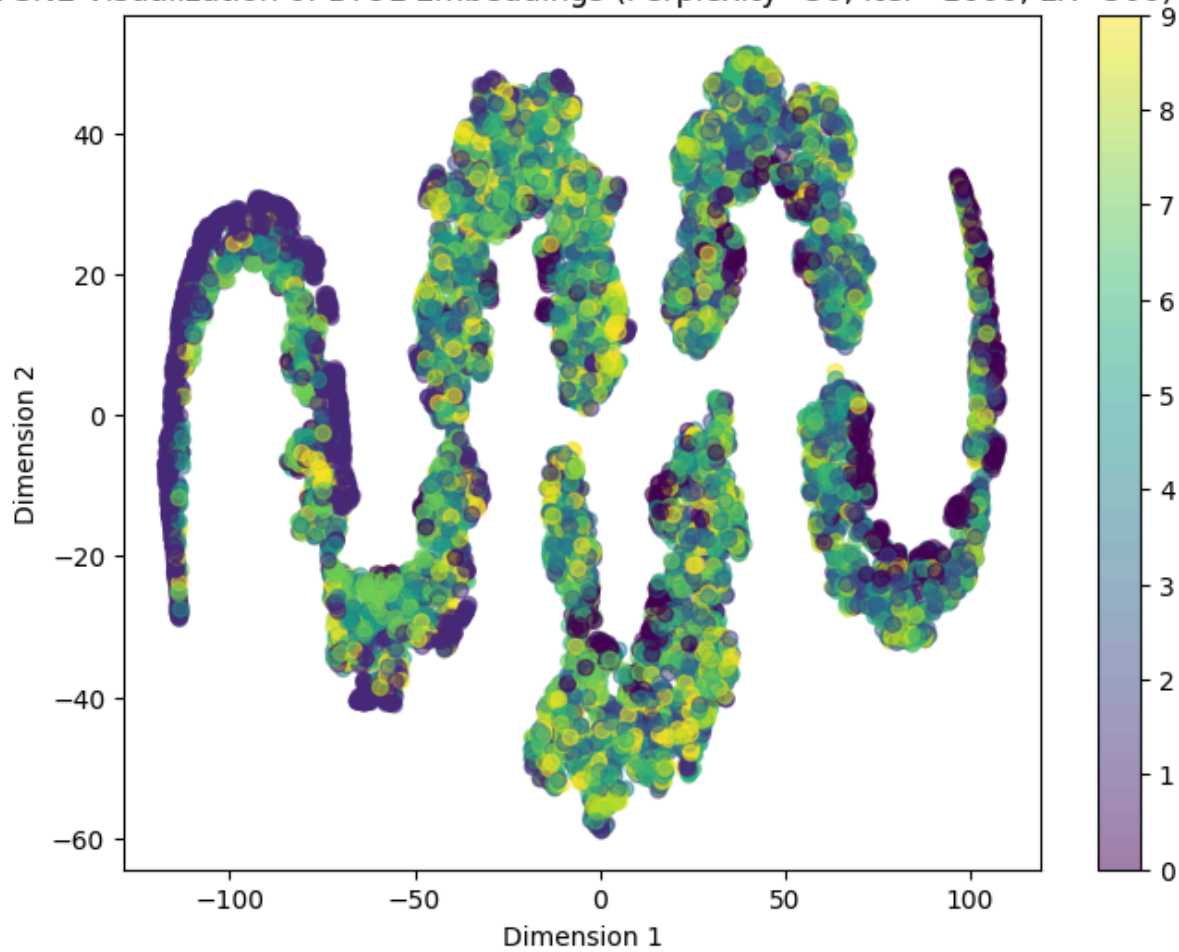
```
C:\Users\13055\AppData\Local\Programs\Python\Python312\Lib\site-packages\sklearn\manifold\_t_sne.py:1162: FutureWarning: 'n_iter' was renamed to 'max_iter' in version 1.5 and will be removed in 1.7.  
  warnings.warn(
```

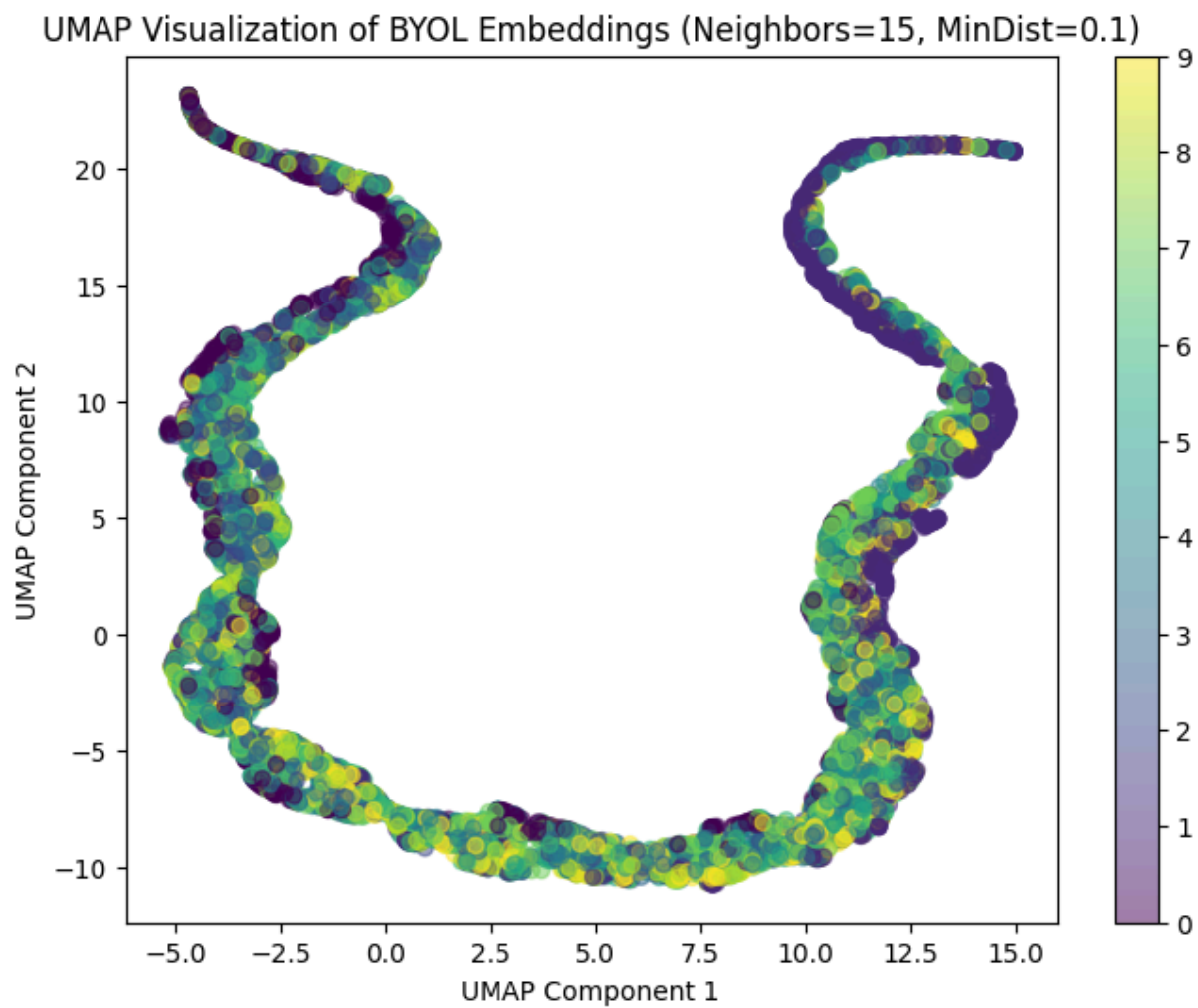
t-SNE Visualization of BYOL Embeddings (Perplexity=30, Iter=1000, LR=200)

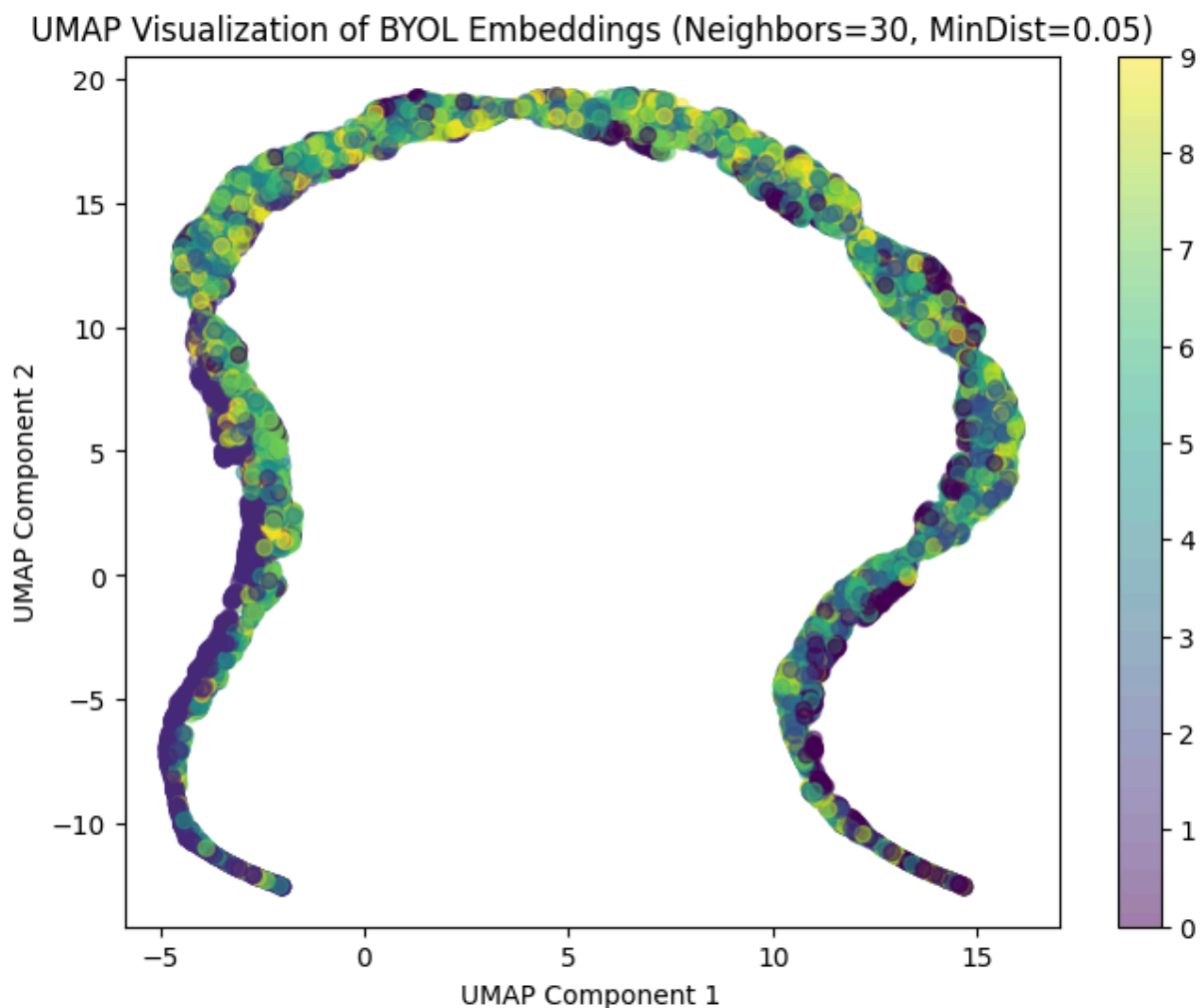


```
C:\Users\13055\AppData\Local\Programs\Python\Python312\Lib\site-packages\sklearn\manifold\_t_sne.py:1162: FutureWarning: 'n_iter' was renamed to 'max_iter' in version 1.5 and will be removed in 1.7.  
  warnings.warn(
```


t-SNE Visualization of BYOL Embeddings (Perplexity=50, Iter=1000, LR=500)







```
In [46]: # Print clustering evaluation results
for method, result in clustering_results.items():
    print(f"{method} Clustering Silhouette Score: {result['score']:.4f}")

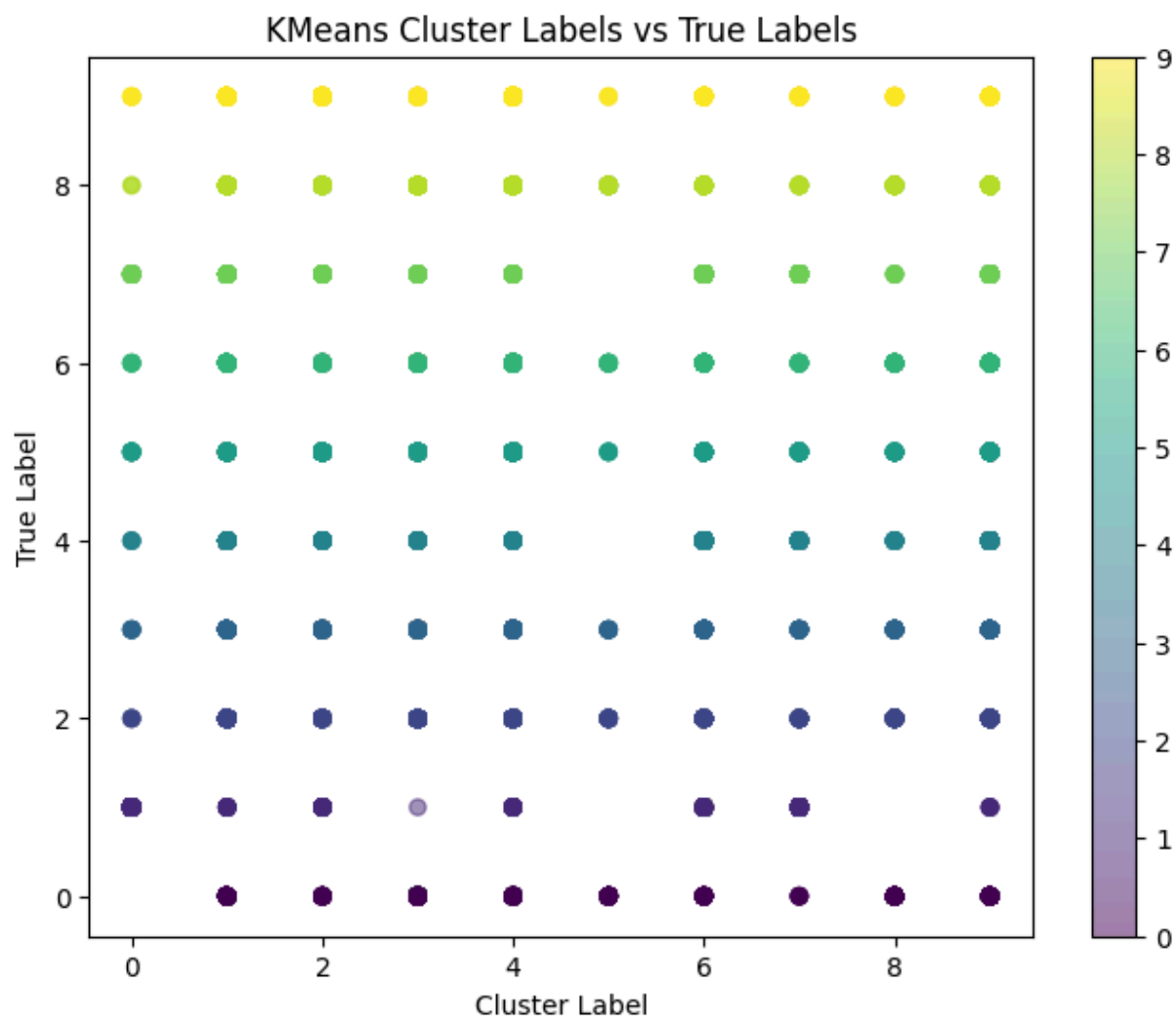
# Compare cluster labels with true labels visually for each clustering method
def visualize_cluster_vs_true_labels(cluster_labels, true_labels, method_name):
    plt.figure(figsize=(8, 6))
    scatter = plt.scatter(cluster_labels, true_labels, c=true_labels, cmap='viridis')
    plt.colorbar(scatter)
    plt.title(f'{method_name} Cluster Labels vs True Labels')
    plt.xlabel('Cluster Label')
    plt.ylabel('True Label')
    plt.show()

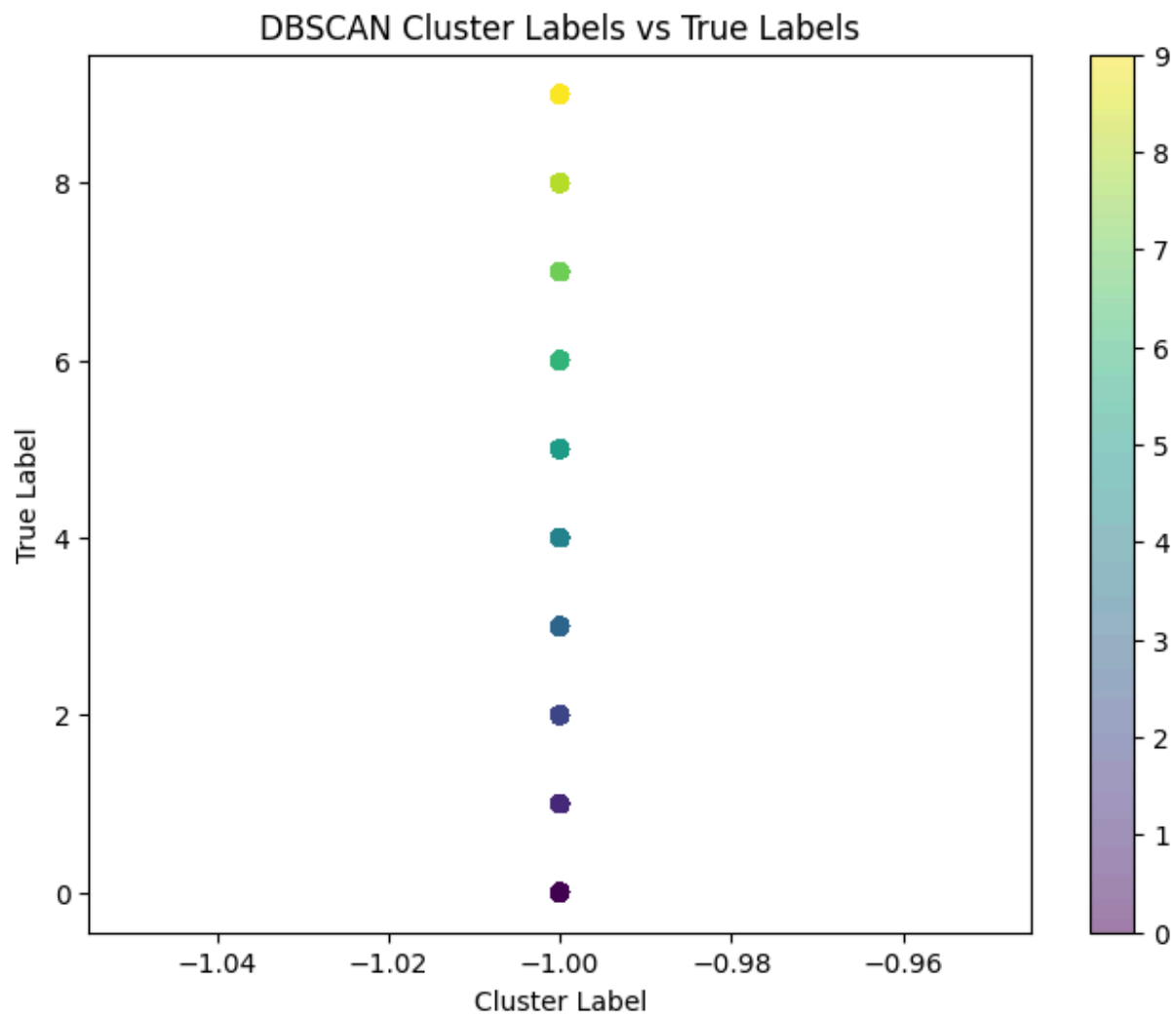
# Visualize clustering results
for method, result in clustering_results.items():
    visualize_cluster_vs_true_labels(result['labels'], test_labels, method)
```

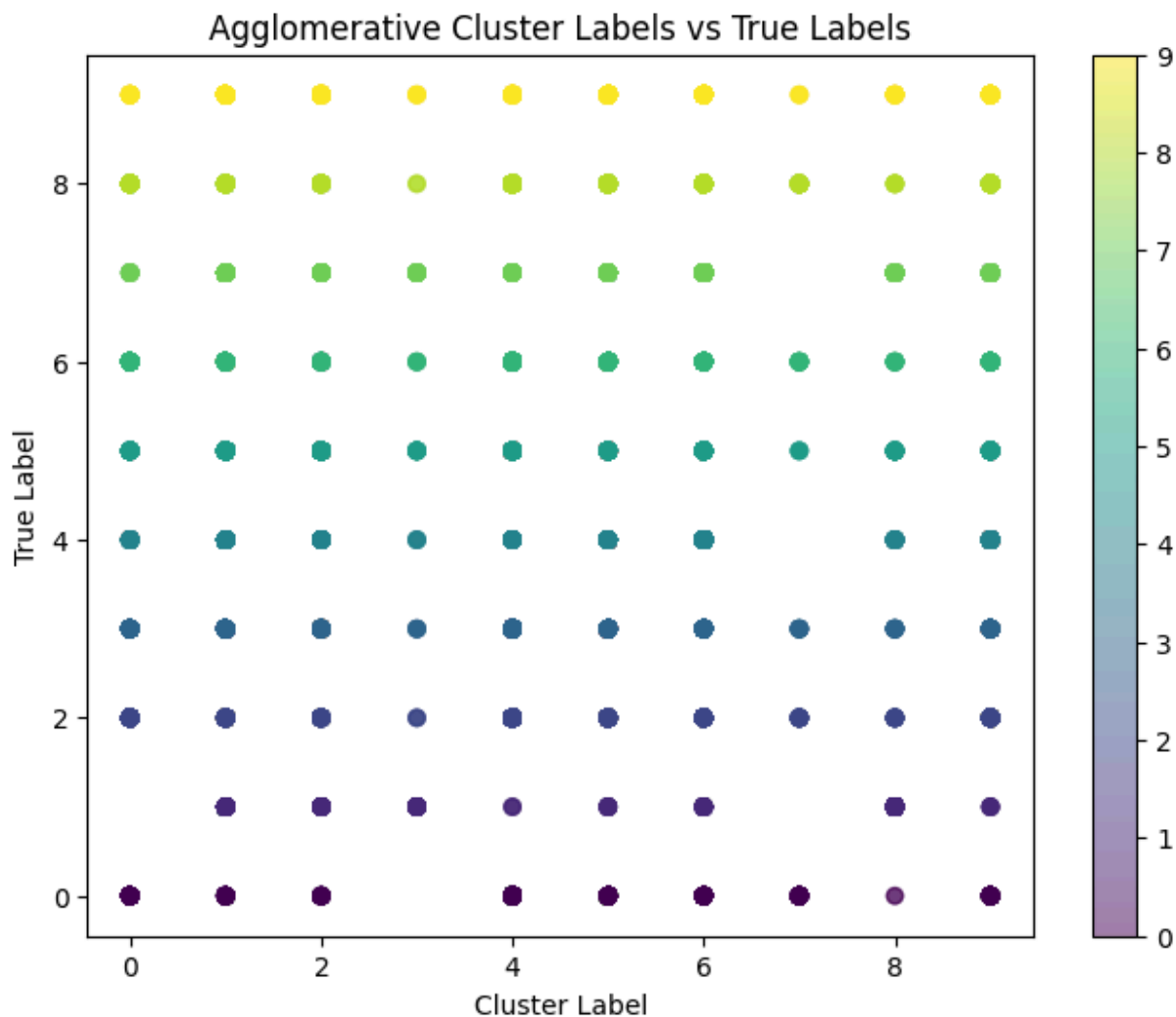
KMeans Clustering Silhouette Score: 0.4727

DBSCAN Clustering Silhouette Score: -1.0000

Agglomerative Clustering Silhouette Score: 0.4487







```
In [47]: # Further Enhanced Fine-Tuning Code
# Fine-tuning is performed to adapt the pre-trained model's Learned features to the
# This allows the model to make minor adjustments to high-level representations bas
# improving performance on the specific classification task without needing to retr

# Step 1: Adjust the model architecture and settings for better fine-tuning

# Freeze all layers of the base model to retain the Learned features from previous
for layer in base_model.layers:
    layer.trainable = False # Freeze base model weights

# Unfreeze the last few layers of the base model to allow fine-tuning
# By freezing lower layers and only fine-tuning the top layers, we reduce the risk
# the general features learned during the initial unsupervised training.
for layer in base_model.layers[-5:]: # Unfreeze the last 5 layers of the base mode
    layer.trainable = True # Allow the weights of these layers to be updated durin

# Modify the classification head with improved architecture
classifier = models.Sequential([
    base_model,
    layers.Flatten(), # Flatten the output of the base model
```

```

layers.Dense(128, kernel_initializer='he_normal'), # He initialization
layers.BatchNormalization(), # Batch normalization before activation
layers.PReLU(), # Parametric ReLU activation for non-linear
layers.Dropout(0.4), # Slightly reduced dropout for regularization
layers.Dense(64, kernel_initializer='he_normal'), # Reduced units
layers.BatchNormalization(), # Batch normalization before activation
layers.PReLU(), # Parametric ReLU activation
layers.Dropout(0.4), # Reduced Dropout
layers.Dense(10, activation='softmax') # Output layer for 10 digit classes
])

# Step 2: Compile the model with adjusted optimizer and learning rate
classifier.compile(optimizer=optimizers.Adam(learning_rate=5e-5), # Even smaller learning rate
                  loss='sparse_categorical_crossentropy',
                  metrics=['accuracy'])

# Adjust Learning rate scheduler to prevent reducing too fast
lr_scheduler = callbacks.ReduceLROnPlateau(monitor='val_loss', factor=0.7, patience=10)

# Early stopping to prevent overfitting with more patience
early_stopping = callbacks.EarlyStopping(monitor='val_loss', patience=6, restore_best_weights=True)

# Input normalization layer to improve model learning
normalization_layer = layers.Normalization(axis=-1)
normalization_layer.adapt(train_images)

# Step 3: Create a new classifier with input normalization
classifier = models.Sequential([
    normalization_layer, # Add normalization as the first layer
    base_model,
    layers.Flatten(),
    layers.Dense(128, kernel_initializer='he_normal'),
    layers.BatchNormalization(),
    layers.PReLU(),
    layers.Dropout(0.4),
    layers.Dense(64, kernel_initializer='he_normal'),
    layers.BatchNormalization(),
    layers.PReLU(),
    layers.Dropout(0.4),
    layers.Dense(10, activation='softmax')
])


# Step 4: Compile the new model
classifier.compile(optimizer=optimizers.Adam(learning_rate=5e-5), # Even smaller learning rate
                  loss='sparse_categorical_crossentropy',
                  metrics=['accuracy'])

# Step 5: Train the model on the Labeled MNIST data for fine-tuning
print("Fine-tuning the model with further enhancements:")
history = classifier.fit(train_images, train_labels, epochs=15,
                        validation_data=(test_images, test_labels),
                        callbacks=[lr_scheduler, early_stopping])


# Step 6: Evaluate the fine-tuned model on the test set
test_loss, test_accuracy = classifier.evaluate(test_images, test_labels)
print(f'Test Accuracy after BYOL and further enhanced fine-tuning: {test_accuracy}')
```

Fine-tuning the model with further enhancements:


Epoch 1/15

1875/1875  20s 4ms/step - accuracy: 0.2002 - loss: 2.4429 - val_accuracy: 0.4648 - val_loss: 1.5584 - learning_rate: 5.0000e-05


Epoch 2/15

1875/1875  7s 4ms/step - accuracy: 0.5006 - loss: 1.4849 - val_accuracy: 0.7902 - val_loss: 0.7895 - learning_rate: 5.0000e-05


Epoch 3/15

1875/1875  7s 4ms/step - accuracy: 0.7156 - loss: 0.9279 - val_accuracy: 0.8826 - val_loss: 0.4882 - learning_rate: 5.0000e-05


Epoch 4/15

1875/1875  8s 4ms/step - accuracy: 0.8168 - loss: 0.6371 - val_accuracy: 0.9130 - val_loss: 0.3737 - learning_rate: 5.0000e-05


Epoch 5/15

1875/1875  7s 4ms/step - accuracy: 0.8652 - loss: 0.4781 - val_accuracy: 0.9333 - val_loss: 0.2597 - learning_rate: 5.0000e-05


Epoch 6/15

1875/1875  7s 4ms/step - accuracy: 0.8941 - loss: 0.3775 - val_accuracy: 0.9305 - val_loss: 0.2381 - learning_rate: 5.0000e-05


Epoch 7/15

1875/1875  7s 4ms/step - accuracy: 0.9124 - loss: 0.3101 - val_accuracy: 0.9538 - val_loss: 0.1717 - learning_rate: 5.0000e-05


Epoch 8/15

1875/1875  7s 4ms/step - accuracy: 0.9247 - loss: 0.2688 - val_accuracy: 0.9541 - val_loss: 0.1545 - learning_rate: 5.0000e-05


Epoch 9/15

1875/1875  8s 4ms/step - accuracy: 0.9332 - loss: 0.2352 - val_accuracy: 0.9562 - val_loss: 0.1469 - learning_rate: 5.0000e-05


Epoch 10/15

1875/1875  8s 4ms/step - accuracy: 0.9394 - loss: 0.2151 - val_accuracy: 0.9592 - val_loss: 0.1365 - learning_rate: 5.0000e-05


Epoch 11/15

1875/1875  8s 4ms/step - accuracy: 0.9420 - loss: 0.2035 - val_accuracy: 0.9586 - val_loss: 0.1409 - learning_rate: 5.0000e-05


Epoch 12/15

1875/1875  8s 4ms/step - accuracy: 0.9465 - loss: 0.1864 - val_accuracy: 0.9751 - val_loss: 0.0894 - learning_rate: 5.0000e-05


Epoch 13/15

1875/1875  7s 4ms/step - accuracy: 0.9498 - loss: 0.1733 - val_accuracy: 0.9667 - val_loss: 0.1125 - learning_rate: 5.0000e-05


Epoch 14/15

1856/1875  0s 3ms/step - accuracy: 0.9526 - loss: 0.1634

Epoch 14: ReduceLROnPlateau reducing learning rate to 3.499999911582563e-05.

1875/1875  7s 4ms/step - accuracy: 0.9526 - loss: 0.1634 - val_accuracy: 0.9637 - val_loss: 0.1213 - learning_rate: 5.0000e-05

Epoch 15/15

1875/1875  7s 4ms/step - accuracy: 0.9560 - loss: 0.1499 - val_accuracy: 0.9646 - val_loss: 0.1151 - learning_rate: 3.5000e-05

Restoring model weights from the end of the best epoch: 12.

313/313  1s 2ms/step - accuracy: 0.9692 - loss: 0.1041

Test Accuracy after BYOL and further enhanced fine-tuning: 97.51%

In [48]: *# This code defines a function to visualize the predictions of a trained neural net
It displays a specified number of images along with their predicted and true labels
the model's performance on individual test cases. By comparing predictions with a
areas where the model performs well and where it struggles, offering insights into*

*# the function calculates the accuracy on the displayed images, giving a quick quan
on the selected subset.*

```
def visualize_predictions(model, images, labels, num_images=20):
    predictions = model.predict(images[:num_images])
    plt.figure(figsize=(15, 3))
    correct_count = 0
    for i in range(num_images):
        predicted_label = np.argmax(predictions[i])
        true_label = labels[i]
        plt.subplot(1, num_images, i + 1)
        plt.imshow(images[i].reshape(28, 28), cmap='gray')
        plt.title(f'Pred: {predicted_label}\nTrue: {true_label}')
        plt.axis('off')
        if predicted_label == true_label:
            correct_count += 1
    plt.show()
    print(f'Correct Predictions: {correct_count}/{num_images}')
    print(f'Accuracy on displayed images: {(correct_count / num_images) * 100:.2f}%')
```

Visualize predictions on test images

```
visualize_predictions(classifier, test_images, test_labels)
```

1/1 ————— 0s 338ms/step

Pred: 7Pred: 2Pred: 1Pred: 0Pred: 4Pred: 1Pred: 4Pred: 9Pred: 5Pred: 9Pred: 0Pred: 6Pred: 9Pred: 0Pred: 1Pred: 5Pred: 9Pred: 7Pred: 8Pred: 4
True: 7True: 2True: 1True: 0True: 4True: 1True: 4True: 9True: 5True: 9True: 0True: 6True: 9True: 0True: 1True: 5True: 9True: 7True: 3True: 4



Correct Predictions: 19/20

Accuracy on displayed images: 95.00%

In []: