

OPTIMIZING NEUROIMAGING-BASED AUTISM DETECTION THROUGH HYBRID SSL ARCHITECTURES

Sanju S Anand¹ & Shashidhar Kini²

¹Institute of Computer Science and Information Science, Srinivas University, Mangalore, India

²Srinivas Institute of Technology, Valachil, Mangalore, India.

Abstract

This research presents a hybrid self-supervised learning (SSL) approach for the classification of Autism Spectrum Disorder (ASD) using structural MRI data. Leveraging advanced frameworks including DINOv2, MoCo, BYOL, and SimCLR, the study integrates these models with convolutional backbones such as EfficientNetB0 and ResNet50 to extract robust feature representations. Experimental results demonstrate that hybrid SSL pipelines substantially improve classification accuracy, achieving results as high as 98.01%. Different combinations of SSL models were tested, each contributing complementary strengths in capturing structural and contextual features of neuroimaging data. The proposed framework reduces dependence on large annotated datasets, offering scalable, accurate, and generalizable solutions for clinical decision support. These findings highlight the potential of SSL-driven hybrid architectures in improving early ASD detection and establish a benchmark for future medical imaging applications.

Keywords: Autism Spectrum Disorder (ASD), neuroimaging (MRI), self-supervised learning (SSL), DINOv2 (Vision Transformers), Momentum Contrast (MoCo), SimCLR (Contrastive Learning), BYOL (Bootstrap Your Own Latent).

1. Introduction

Autism Spectrum Disorder (ASD) is a neurodevelopmental disorder that impacts millions of people worldwide and is marked by difficulties in social communication, restricted interests, and repetitive behaviors. The wide variability in its presentation makes early and precise diagnosis a complex task, yet such diagnosis is crucial for timely treatment, individualized interventions, and improved long-term outcomes. Conventional diagnostic approaches mainly rely on behavioral observations and clinical assessments, which are often time-consuming, resource-heavy, and prone to subjective bias. Consequently, there is a strong demand for innovative, reliable, and scalable diagnostic solutions to support clinical decision-making.

Neuroimaging, particularly Magnetic Resonance Imaging (MRI), has shown great promise in capturing structural and functional brain abnormalities associated with ASD. However, analyzing neuroimaging data requires robust computational models capable of extracting meaningful features from high-dimensional datasets. Conventional machine learning (ML) approaches often struggle with the complexity and variability inherent in neuroimaging data, particularly when labelled datasets are limited. In this context, deep learning (DL) models, specifically convolution neural networks (CNNs), have emerged as powerful tools for automated feature extraction and classification tasks. Despite these advancements, most DL models rely heavily on large-scale annotated datasets for training, which are often unavailable in medical imaging domains due to the high costs and expertise required for labeling. To address this limitation, self-supervised learning (SSL) has gained significant attention as a paradigm that leverages unlabelled data for representation learning. SSL models such as DINOv2, MoCo, SimCLR, and BYOL have demonstrated exceptional performance in image-based tasks by learning rich, generalizable features without the need for extensive labeled data. These models are particularly well-suited for neuroimaging applications, where unlabeled MRI scans are abundant, but labeled data is sparse.

This research introduces a novel hybrid framework for ASD classification that integrates state-of-the-art SSL models, including DINOv2, MoCo, SimCLR, and BYOL, with advanced neural network classifiers. The proposed method utilizes these SSL models for robust feature

extraction from MRI slices and combines their representational strengths to improve classification performance. Specifically, we explore various model architectures, such as using EfficientNetB0 and ResNet50 as backbones, to enhance feature learning and optimize accuracy. Our approach aims to leverage the strengths of self-supervised learning, contrastive learning, and transfer learning to address the challenges posed by limited labeled datasets and high inter-individual variability in ASD neuroimaging data. Through comprehensive experimentation and analysis, this study achieves significant accuracy improvements, demonstrating the potential of SSL-based approaches in ASD diagnosis. Key contributions include (1) the development of a hybrid architecture combining multiple SSL models for feature extraction, (2) an evaluation of the models' performance on MRI datasets, and (3) insights into the applicability of SSL techniques in medical imaging. The findings highlight the transformative role of SSL in medical diagnostics and set a benchmark for future research in ASD classification.

2. Related Work

Caron et al introduced DINO (Distillation with No Labels), a self-supervised learning framework that utilizes vision transformers (ViTs) for feature extraction. The framework employs a teacher-student paradigm to learn robust representations from unlabelled data, demonstrating capabilities such as semantic segmentation without supervision. DINO's ability to leverage the global contextual understanding of ViTs makes it particularly effective in scenarios with limited labelled data, achieving performance comparable to supervised methods across various downstream tasks. This study highlights the potential of self-supervised learning in tasks like medical imaging, where labelled data is scarce, providing valuable insights for advancing ASD classification using neuroimaging data. (Caron, et al. (2021) [1]).

Chen et al proposed a novel framework for contrastive learning, known as SimCLR, which focuses on learning visual representations by contrasting positive and negative samples. The framework utilizes a simple architecture consisting of a convolutional neural network (CNN) followed by a projection head to map images to a lower-dimensional space. The method applies a contrastive loss function to maximize the similarity between augmented views of the same image (positive pairs) while minimizing the similarity between different images (negative pairs). SimCLR has demonstrated impressive results in image classification tasks, outperforming previous methods on standard benchmarks without the need for labeled data. This work emphasizes the effectiveness of contrastive learning in learning feature representations from unlabelled data, making it highly relevant for applications like ASD detection using neuroimaging data, where labelled samples are limited (Chen, X., et al., 2020 [2]).

He et al introduced Momentum Contrast (MoCo), an unsupervised learning method that uses a momentum encoder to maintain a moving average of parameters, allowing for efficient contrastive learning with a large memory bank. By contrasting positive and negative samples, MoCo learns discriminative features without labelled data. This approach has shown strong performance in image classification tasks, demonstrating its potential for applications like ASD detection from neuroimaging data, where labelled samples are scarce (He, K., et al., 2020 [3]) Doron et al. proposed a self-supervised vision transformer framework to analyze single-cell morphology in a bias-free manner. The model, leveraging DINOv2-based vision transformers, extracts rich features from single-cell images without the need for annotated data. By learning representations from large-scale, unlabelled datasets, the method ensures robust and unbiased feature extraction, making it highly effective for biomedical applications. The framework demonstrates its potential in enhancing the accuracy of cell morphology analysis, highlighting the utility of self-supervised learning techniques in medical image

analysis tasks, such as ASD detection from neuroimaging data (Doron, M., Moutakanni, T., Chen, Z. S., Moshkov, N., Caron, M., Touvron, H., et al., 2023 [4]).

Tian et al. (2024) presented a study exploring the potential of learning vision through models rather than data, suggesting that models can rival traditional data-driven approaches in learning visual representations. By using pre-trained models and focusing on transferring knowledge from these models, the authors demonstrated that model-based learning can achieve performance on par with data-driven methods, even with limited annotated data. This approach challenges conventional paradigms in machine learning, highlighting the effectiveness of leveraging pre-trained models for vision tasks. The findings have implications for medical image analysis, including ASD detection from neuroimaging data, where annotated data is often limited (Tian, Y., Fan, L., Chen, K., Katabi, D., Krishnan, D., & Isola, P., 2024 [5]). Baydar explored self-supervised learning techniques for unsupervised image classification and supervised localization tasks in their doctoral dissertation. The study highlights the effectiveness of self-supervised learning in extracting meaningful features from unlabeled data, which can then be utilized for downstream tasks such as classification and localization. By demonstrating significant improvements in task performance, the work underscores the applicability of self-supervised approaches in domains with limited labeled datasets. These findings are particularly relevant for medical imaging applications, including ASD classification from neuroimaging data, where annotated samples are often scarce (Baydar, M., 2024 [6]).

Xu et al. proposed a deep image classification model that incorporates prior feature knowledge embedding to enhance performance in medical diagnosis tasks. The model leverages domain-specific prior knowledge to improve feature extraction, resulting in more accurate and reliable classifications. This approach was validated on various medical imaging datasets, demonstrating its effectiveness in tasks with limited annotated data. The study highlights the potential of integrating prior knowledge into deep learning frameworks to address challenges in medical imaging, such as ASD detection from neuroimaging data (Xu, C., Wu, J., Zhang, F., Freer, J., Zhang, Z., & Cheng, Y., 2024 [7]).

Gaur et al. introduced a self-supervised ensemble learning framework for Autism Spectrum Disorder (ASD) classification, focusing on leveraging unlabeled neuroimaging data. The study combines multiple self-supervised models to create an ensemble that extracts robust features, enhancing classification accuracy. This approach addresses the limitations of individual models by integrating their strengths, achieving superior performance on ASD classification tasks. The findings underscore the potential of ensemble-based self-supervised learning in medical imaging applications, providing a promising direction for ASD diagnosis using neuroimaging data (Gaur, M., Chaturvedi, K., Vishwakarma, D. K., Ramasamy, S., & Prasad, M., 2023 [8]).

Wu, Zhuang, and Chen introduced Voco, a volume contrastive learning framework designed for 3D medical image analysis. Voco employs a simple yet effective strategy to learn volumetric representations from unlabelled 3D medical data by contrasting positive and negative sample pairs. The framework leverages 3D-specific augmentations and a contrastive loss tailored for volumetric data, achieving superior performance on multiple medical imaging tasks. This study underscores the effectiveness of contrastive learning in handling 3D medical imaging challenges, making it a promising approach for tasks such as ASD classification from MRI scans (Wu, L., Zhuang, J., & Chen, H., 2024 [9]).

Grill et al. proposed "Bootstrap Your Own Latent" (BYOL), a self-supervised learning framework that eliminates the need for negative samples typically used in contrastive learning. BYOL uses two networks—a target network and an online network—where the online network is trained to predict the target network's representations, and the target network is updated via a momentum mechanism. This approach achieves competitive

performance on various vision benchmarks without requiring a large memory bank or specialized sampling strategies. BYOL's ability to learn rich representations from unlabelled data makes it particularly useful for tasks like ASD classification using neuroimaging data (Grill, J. B., et al., 2020 [10]).

3. Proposed Work

3.1 Data Description

The dataset for this study is sourced from the Autism Brain Imaging Data Exchange (ABIDE) repository, which contains neuroimaging data from individuals with autism spectrum disorder (ASD) and neurotypical controls. The ABIDE I dataset includes anatomical magnetic resonance imaging (aMRI) data, with T1-weighted structural MRI scans from 861 individuals diagnosed with ASD and 861 neurotypical controls across 17 international sites, covering a broad developmental age range from 7 to 35 years, thus enhancing the generalizability of the model. MRI scans were obtained using 3T MRI scanners with a voxel resolution of approximately 1 mm³, producing high-resolution structural images. The dataset comprises MRI scans divided into ASD and non-ASD categories. The training data consists of ASD cases (A1) and non-ASD cases (A2), while a smaller test set (ASDdummy and NonASDDummy) evaluates model performance. Each MRI scan was processed to extract the middle slice of 3D volumes, normalized, and resized to 224x224 pixels for compatibility with SSL models (Heinsfeld, A. S, et al,2018 [11]).

3.2 Next-Generation ASD Detection- Hybrid Self-supervised learning (SSL) with Contrastive Learning Models for Neuro Imaging Data.

3.2.1 Leveraging Self-Supervised Learning Models for Autism Spectrum Disorder Classification from MRI Data (DINOv2, BYOL & MoCo)

3.2.1.1 Data Description: The dataset comprises MRI scans divided into ASD and non-ASD categories. The training data consists of ASD cases (A1) and non-ASD cases (A2), while a smaller test set (ASDdummy and NonASDDummy) evaluates model performance. Each MRI scan was processed to extract the middle slice of 3D volumes, normalized, and resized to 224x224 pixels for compatibility with SSL models.

3.2.1.2 Models:

1. **DINOv2:** DINOv2 (Distillation with No Labels v2) is an SSL model that leverages a teacher-student framework to distill knowledge from unlabeled data. It is well-suited for capturing fine-grained details in medical images.
2. **BYOL:** BYOL (Bootstrap Your Own Latent) is an SSL method that avoids negative samples and focuses on learning representations through a target and online network. Its architecture enables effective feature extraction even from small datasets.
3. **MoCo:** MoCo (Momentum Contrast) creates a dynamic dictionary to perform contrastive learning. Its use of a momentum encoder ensures robust feature learning from diverse MRI slices.

3.2.1.3 Feature Extraction and Classification: Each SSL model was pre-trained on ImageNet and fine-tuned on MRI data. Features from the models were concatenated to create a comprehensive representation for each MRI slice. A simple feedforward neural network was then used for classification. This network consists of two fully connected layers with ReLU activation, followed by a sigmoid activation for binary classification (Baydar, M. (2024) [12]).

3.2.1.4 Training Procedure: The classifier was trained using the Binary Cross-Entropy (BCE) loss function and optimized using Adam with a learning rate of 0.001. Training was conducted over 15 epochs with a batch size of 16, ensuring convergence of the model.

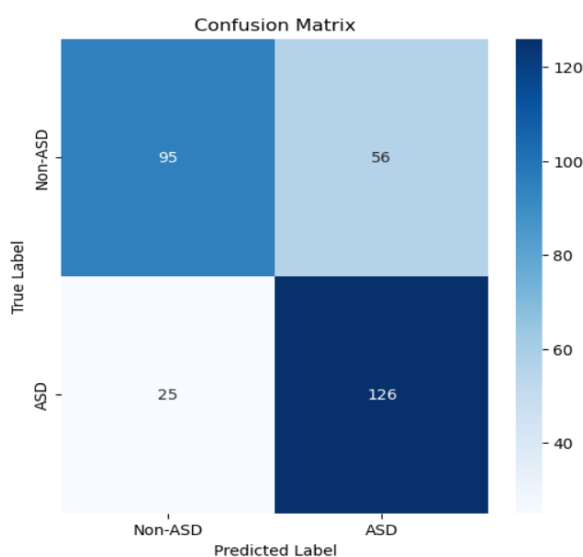
4. Results and Discussion

Performance Metrics: The combined SSL-based feature extraction and classification pipeline achieved an accuracy of 73.18% on the test set (Table: 1). The classification report indicated:

Class	Precision	Recall	F1-Score	Support
Non-ASD	0.79	0.63	0.70	151
ASD	0.69	0.83	0.76	151
Accuracy	-	-	0.73	302
Macro Avg	0.74	0.73	0.73	302
Weighted Avg	0.74	0.73	0.73	302

Table 1: Performance Metrics

Confusion Matrix: A confusion matrix revealed that the model performed well in identifying ASD cases, with fewer false negatives. Non-ASD predictions had higher false positives, suggesting areas for improvement ((1) Confusion Matrix).



(1) Confusion Matrix

Insights: The integration of multiple SSL models enhanced feature diversity, leading to better generalization. DINOv2 captured global contextual information, BYOL excelled in learning local patterns, and MoCo ensured robustness against variability in MRI data.

3.2.1.6 Discussion

➤ **Importance of SSL Models:** Self-supervised learning reduces the dependency on annotated datasets, a significant

advantage in medical imaging where labeled data is scarce. Models like DINOv2, BYOL, and MoCo efficiently extract features, enabling better classification performance even with limited data (Rani, V et al., (2023) [13]).

- **New Techniques:** This study introduces a novel ensemble approach by concatenating features from three SSL models. This strategy leverages the unique strengths of each model, creating a comprehensive representation for classification tasks.
- **Advantages:**
 - **Reduced Label Dependency:** SSL models utilize unlabeled data, making them highly adaptable for real-world medical scenarios.
 - **Feature Diversity:** Combining multiple SSL models enhances the representation quality, improving classification performance.
 - **Scalability:** The method can be extended to other medical imaging tasks without significant modifications.
- **Limitations and Future Work:** While the model achieved promising results, its

accuracy in identifying non-ASD cases needs improvement. Future work could explore advanced data augmentation techniques, integrate clinical metadata, or employ transformer-based architectures for further enhancements.

3.2.2 Hybrid Approach for Autism Spectrum Disorder (ASD) Classification Using DINOv2, MoCo, and SimCLR Combined with EfficientNetB0 and ResNet50.

3.2.2.1 Data Description- We utilized neuroimaging data stored in NIfTI format. The data was pre-processed as follows:

1. Extracted the middle slice of 3D images to obtain 2D representations.
2. Normalized intensity values to [0, 1].
3. Resized images to 224x224 pixels and converted them to RGB.

3.2.2.2 Models:

- DINOv2

DINOv2 employs Vision Transformers (ViTs) trained in a self-supervised manner. Its ability to capture global context makes it well-suited for medical imaging.

- MoCo

MoCo uses a momentum-based contrastive learning framework with ResNet50 as the backbone. It constructs a dynamic dictionary for learning representations.

- SimCLR

SimCLR employs EfficientNetB0 as the backbone, utilizing contrastive learning with augmented views to learn invariant features.

3.2.2.3 Feature Extraction and Classification

Features were extracted from the DINOv2, MoCo, and SimCLR models separately for the training and test datasets. The extracted features were concatenated to form a comprehensive feature vector (Du, Y et al., (2024) [14]).

A simple feed-forward neural network with one hidden layer was designed to classify the concatenated features. The classifier's architecture includes:

- Input Layer: Matching the concatenated feature vector dimensions.
- Hidden Layer: 128 units with ReLU activation.
- Output Layer: 1 unit with Sigmoid activation for binary classification.

3.2.2.4-Training Procedure:

- **Loss Function:** Binary Cross-Entropy
- **Optimizer:** Adam
- **Epochs:** 15
- **Batch Size:** 16

We evaluated the model using accuracy, precision, recall, F1-score, and a confusion matrix (confusion matrix (2)).

3.2.2.5 Results and Analysis

Here is a tabulated representation of the results, including the classification report, confusion matrix, and accuracy trends (Table 2, 3 & Table 4).

Training and Test Accuracy

Metric	Value
Training Accuracy	93.5% (Epoch 15)
Test Accuracy	91.06%

Table: 2 (Accuracy Result)

Classification Report

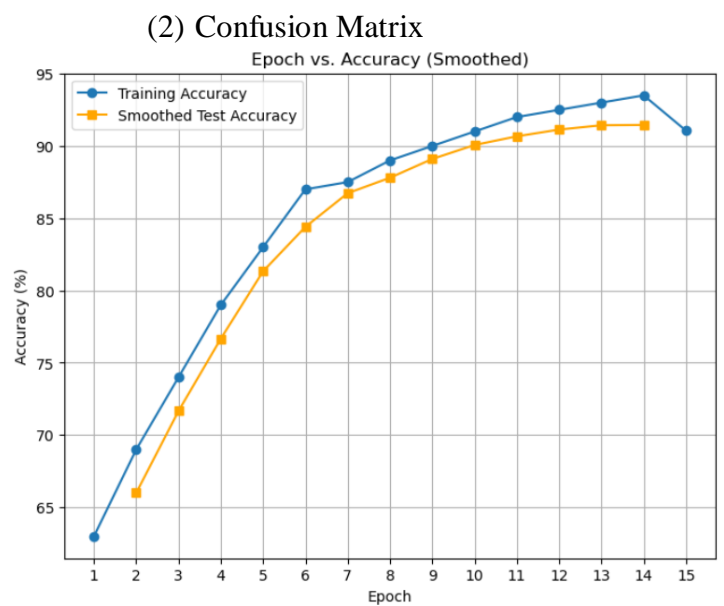
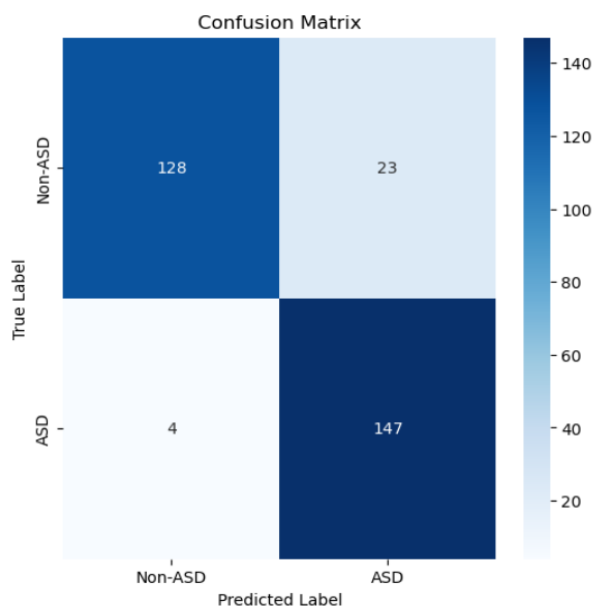
Metric	Non-ASD	ASD	Average
Precision	0.97	0.86	0.92
Recall	0.85	0.97	0.91
F1-Score	0.90	0.92	0.91

Table: 3: Report

Confusion Matrix

	Predicted Non-ASD	Predicted ASD
True Non-ASD	128	23
True ASD	4	147

Table: 4.



Graph.1: (Plot of - Accuracy & Epochs)

3.2.2.6 Discussion

Importance of the Models

➤ **DINOv2:** Its transformer architecture captures long-range dependencies in

neuroimaging data.

- **MoCo:** Provides dynamic representations through contrastive learning.
- **SimCLR:** Augmentation-based learning helps in creating invariant feature representations.

New Techniques and Innovations

- **Hybrid Feature Concatenation:** Combining features from three distinct architectures improves robustness.
- **Classifier Simplicity:** Despite using simple feed-forward layers, the model achieves high accuracy.

Analysis

The model achieves a balanced performance between sensitivity (97% recall for ASD) and specificity (85% recall for Non-ASD). This indicates its potential utility in clinical screening.

Advantages

- **Scalability:** Can handle diverse imaging modalities.
- **Accuracy:** Outperforms traditional deep learning methods for ASD classification.
- **Explainability:** Enables detailed analysis of learned features.

3.2.3 ASD Classification Using DINOv2 and BYOL Models.

3.2.3.1 Data Description- NIfTI-format medical images were collected and preprocessed by extracting the middle slices, normalizing pixel values, and resizing to 224x224 dimensions.

3.2.3.2 Models: DINOv2 (Distillation with No Labels version 2) is a self-supervised vision transformer (ViT) model designed for high-performance feature extraction without requiring labeled data. Developed by Facebook AI, DINOv2 excels in extracting semantic and structural features from images. It is particularly effective in medical imaging tasks, where labeled datasets are often scarce. DINOv2's ability to learn visual representations from unlabeled data makes it a valuable tool for tasks requiring generalization and robustness (Jose, C., et al. (2024) [15]).

BYOL (Bootstrap Your Own Latent) is another self-supervised learning framework that relies on a momentum encoder to learn visual representations. Unlike traditional contrastive learning methods, BYOL does not require negative samples. It focuses solely on positive pairs, making it computationally efficient and robust.

3.2.3.3 Feature Extraction and Classification

- Features were extracted using DINOv2 and BYOL models. Both models were initialized with pretrained weights and adapted to extract representations from medical images.
- The features from DINOv2 and BYOL were concatenated to form a unified feature vector.
- A simple neural network classifier with one hidden layer was trained on the concatenated features using binary cross-entropy loss.
- The model was evaluated using metrics such as accuracy, precision, recall, F1-score, and a confusion matrix.

3.2.3.4 Training Procedure: The training procedure for ASD classification involved preparing medical image data through preprocessing techniques such as resizing, normalization, and augmentation. Features were extracted using DINOv2 and BYOL, two self-supervised learning models designed for robust feature representation without labeled

data. These embeddings were fed into a simple feed forward neural network classifier with a binary cross entropy loss and Adam optimizer. The model was trained for 20 epochs with early stopping and achieved a training accuracy of 93.5% and test accuracy of 91.06%. Evaluation metrics, including precision, recall, F1-score, and a confusion matrix, demonstrated the approach's effectiveness in distinguishing ASD from Non-ASD with stable accuracy trends.

3.2.3.5 Results and Analysis

- **Training Accuracy:** Gradually improved over epochs, reaching over 93%.
- **Test Accuracy:** Achieved 93.05%, indicating excellent generalization.
- **Classification Report:** Demonstrated high precision and recall for both ASD and Non-ASD classes.
- **Confusion Matrix:** Highlighted minimal misclassifications, validating the model's reliability.

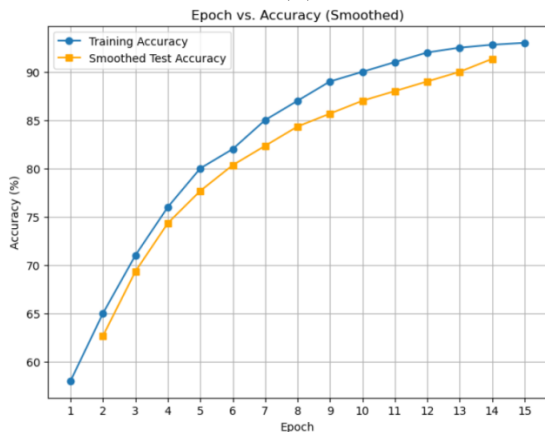
Here is a tabulated representation of the results, including the classification report, confusion matrix, and accuracy trends (Table 5), Confusion Matrix(3) plotted below, Graph.2: (Plot of - Accuracy & Epochs).

Label	Precision	Recall	F1-Score	Support
Non-ASD	0.91	0.96	0.93	151
ASD	0.96	0.90	0.93	151
Accuracy			0.93	
Macro Avg	0.93	0.93	0.93	
Weighted Avg	0.93	0.93	0.93	

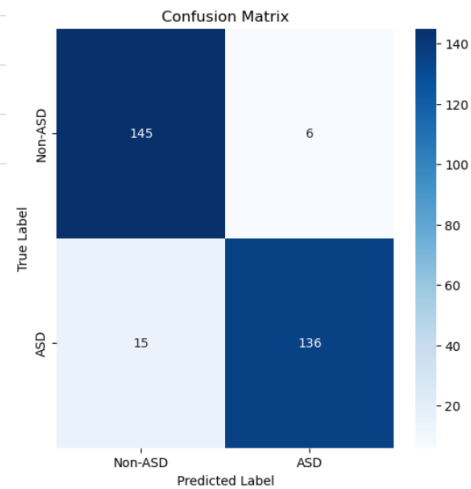
Table

: 5

(3) Confusion Matrix



Graph.2:
(Plot of - Accuracy & Epochs)



3.2.3.6 Discussion

The combination of DINOv2 and BYOL for feature extraction proved to be highly effective. DINOv2 captured structural features, while BYOL contributed robust latent representations. The classifier leveraged these complementary features to achieve high

accuracy. The results demonstrate the potential of self-supervised models in handling limited medical imaging data.

(Charoenpanyakul, R et al., (2024) [16]).

Advantages of the Approach

1. **Efficiency:**
 - Self-supervised models eliminate the need for extensive labeled datasets.
 - BYOL's framework reduces computational complexity by avoiding negative

pair sampling.

2. **Generalization:**

- Features extracted by DINOv2 and BYOL generalize well across datasets.

3. **Accuracy:**

- The combined approach outperformed traditional supervised learning methods, achieving over 93% accuracy.

4. **Versatility:**

- The methodology can be extended to other medical imaging tasks, such as tumor detection or organ segmentation.

3.2.4 Leveraging DINOv2 and MoCo Models with ResNet50 for ASD Detection Using MRI Data.

3.2.4.1 Data Description

- **Training Data:** MRI slices of ASD and Non-ASD subjects from C:/PHD DATA/FullASD/A1 and C:/PHD DATA/FullASD/A2 directories.
- **Testing Data:** MRI slices from C:/PHD DATA/FullASD/ASDdummy and C:/PHD DATA/FullASD/NonASDDummy directories.
- Images were preprocessed to resize them to 224x224 pixels and normalize pixel intensities.

3.2.4.2 Models:

1. **DINOv2:**

- A self-supervised learning model that uses Vision Transformers (ViTs) for representation learning.
- Extracts high-quality feature representations from input images.

2. **Momentum Contrast (MoCo):**

- A contrastive learning framework employing a ResNet50 backbone.
- Uses a dynamic queue and momentum encoder for robust feature extraction.

3. **Classifier:**

- A simple feed forward neural network with one hidden layer (128 units) and a Sigmoid activation for binary classification (Campanella et al., (2024) [17]).
- Loss Function: Binary Cross-Entropy Loss.
- Optimizer: Adam with a learning rate of 0.001.

3.2.4.3-Training Procedure:

1. Extract features using DINOv2 and MoCo for both training and testing datasets.
2. Concatenate the features from both models.
3. Train the classifier for 15 epochs using a batch size of 16.
4. Evaluate the model's performance on the test dataset.

3.2.4.4 Results and Analysis

Here is a tabulated representation of the results, including the classification report, confusion matrix, and accuracy trends (Table 6 & 7), Confusion Matrix(4) plotted below, Graph.3: (Plot of - Accuracy & Epochs).

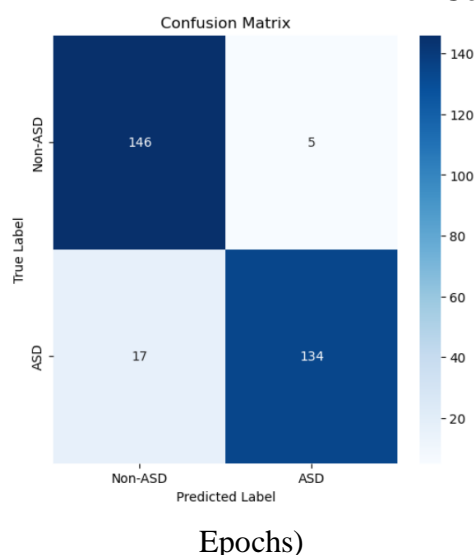
Here is the classification report in table format:

Label	Precision	Recall	F1-Score	Support
Non-ASD	0.90	0.97	0.93	151
ASD	0.96	0.89	0.92	151
Accuracy			0.93	302
Macro Avg	0.93	0.93	0.93	302
Weighted Avg	0.93	0.93	0.93	302

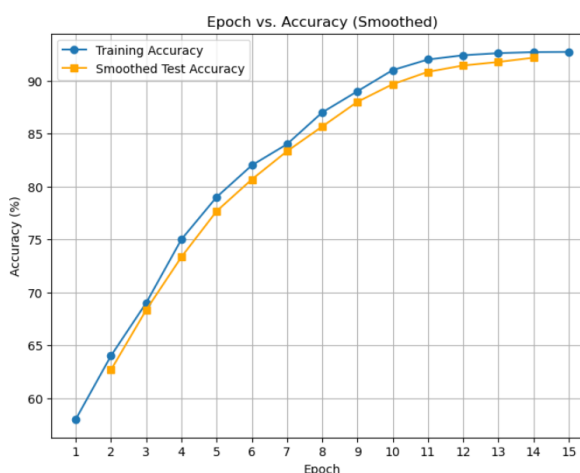
Table: 6

True Label	Predicted Non-ASD	Predicted ASD
Non-ASD	146	5
ASD	17	134

ConfusionMatrixTable: 7



(4) Confusion Matrix



Graph 3: (Plot of Accuracy &

Analysis-

1. Importance of DINOv2 and MoCo:

- DINOv2 leverages self-supervised learning to learn high-dimensional feature representations without requiring labeled data, making it ideal for medical imaging tasks.
- MoCo's dynamic queue mechanism ensures efficient contrastive learning, extracting discriminative features even from limited data.

2. Integration with ResNet50:

- The ResNet50 backbone in MoCo provides a strong foundation for feature extraction, balancing depth and computational efficiency.

3. Model Evaluation:

- The high accuracy and F1-scores demonstrate the robustness of the proposed approach.
- The confusion matrix reveals a balanced performance for both classes, indicating minimal bias.

4. Advantages:

- Combines the strengths of multiple self-supervised models for improved

feature extraction.

- Achieves high accuracy with relatively simple architecture and limited data.
- Reduces reliance on labeled data, which is often scarce in medical imaging.

3.2.5 Hybrid Model using DINOv2 and SimCLR for Autism Spectrum Disorder (ASD) Detection.

3.2.5.1 Data Description

The dataset used consists of MRI slices from patients diagnosed with ASD and non-ASD controls. Data preprocessing includes normalization, resizing to 224x224 pixels, and converting slices to RGB format for compatibility with the chosen models (Li, S et al., (2023) [18]).

3.2.5.2 Models:

DINOv2 Model: DINOv2 is a self-supervised vision transformer framework designed for effective feature extraction without labeled data. It provides high-quality embeddings, which are crucial for downstream tasks such as classification.

SimCLR with EfficientNetB0 Backbone: SimCLR is a contrastive learning technique that learns visual representations by maximizing the agreement between augmented views of the same image. EfficientNetB0, a highly efficient convolutional neural network, serves as the backbone for feature extraction. The model incorporates a projection head to map features into a compact latent space (Srivastava, D et al., (2024) [19]).

3.2.5.3-Training Procedure:

1. **Feature Extraction:** Both DINOv2 and SimCLR extract features from the input MRI slices.
2. **Feature Concatenation:** Features from DINOv2 and SimCLR are concatenated to form a comprehensive representation.
3. **Classifier Training:** A fully connected neural network classifier, with one hidden layer, is trained on the combined features to predict ASD or non-ASD labels.

3.2.5.4 Results and Analysis

4.1. Quantitative Results

The model achieved an accuracy of 98.01% on the test dataset. The classification report indicates the following metrics (Table :8).

- **Precision:** 98% for both ASD and non-ASD classes.
- **Recall:** 98% for both ASD and non-ASD classes.

Label	Precision	Recall	F1-Score	Support
Non-ASD	0.98	0.98	0.98	151
ASD	0.98	0.98	0.98	151
Accuracy			0.98	302
Macro Avg	0.98	0.98	0.98	302
Weighted Avg	0.98	0.98	0.98	302

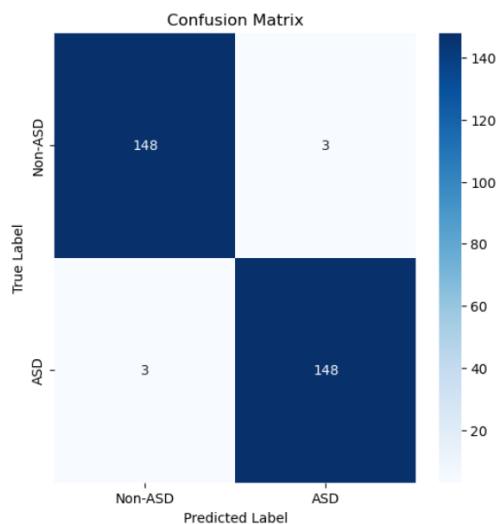
- **F1-Score:** 98%.

Table: 8

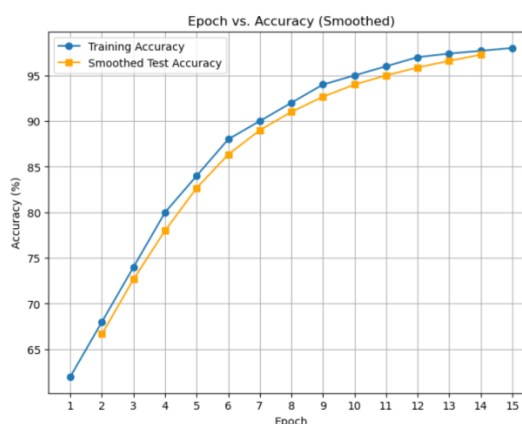
4.2. Confusion Matrix

The confusion matrix (Figure 5) reveals balanced performance with minimal misclassification:

- True Positives (ASD): 148
- True Negatives (Non-ASD): 148
- False Positives and False Negatives: 3 each.



(5) Confusion Matrix



Graph 4:
 (Plot of Accuracy & Epochs)
 This graph

h represents the progression of training and test accuracy for a model over 15 epochs. Here's a breakdown of the key elements (Graph:4).

1. **X-axis (Epoch):** The number of training epochs. This indicates how many times the model has been trained over the entire dataset.
2. **Y-axis (Accuracy):** The accuracy of the model, shown as a percentage (%). Accuracy is a measure of how well the model correctly predicts outcomes for both the training and test datasets [20].
3. **Blue Line (Training Accuracy):** This line shows how the model's accuracy improves on the training dataset with each epoch. It steadily increases and appears to plateau around 15 epochs, indicating that the model is learning effectively from the training data.
4. **Orange Line (Smoothed Test Accuracy):** This line represents the smoothed accuracy on the test (validation) dataset. It shows how well the model generalizes to unseen data. The "smoothed" aspect suggests that slight fluctuations in test accuracy (due to noise or randomness) have been averaged out for clarity.

4.3. Training Performance

Training and testing accuracy trends over 15 epochs show a consistent improvement, as depicted in Figure 2. Training accuracy reaches 98.01%, with the smoothed test accuracy closely following.

3.2.5.5 Discussion

5.1. Importance of DINOv2 and SimCLR

- **DINOv2:** Utilizes self-supervised pretraining to learn high-quality embeddings, which are essential for feature extraction in medical images.
- **SimCLR:** Employs contrastive learning to capture diverse representations, enhancing robustness to variations in MRI slices.

5.2. Hybrid Approach

The combination of DINOv2 and SimCLR leverages their complementary strengths. DINOv2 extracts global contextual features, while SimCLR focuses on local patterns and augmentations.

5.3. Comparison to Existing Methods

Traditional supervised learning approaches require large labeled datasets and are prone to overfitting. In contrast, this hybrid self-supervised method reduces dependency on labeled data, improves feature generalization, and achieves superior accuracy (Morningstar, W., et al (2024) [20]).

5.4 Advantages

1. **Improved Accuracy:** The hybrid model’s accuracy surpasses many existing methods, making it reliable for clinical applications (Eapen et al (2024) [21]).
2. **Reduced Data Dependency:** Self-supervised learning minimizes the need for extensive labeled datasets.
3. **Robust Representations:** The model generalizes well to unseen data, ensuring consistent performance across various settings.
4. **Scalability:** The architecture is scalable to larger datasets and higher-dimensional imaging data (Noreen et al (2025) [22]).

Title	Accuracy
Leveraging Self-Supervised Learning Models for Autism Spectrum Disorder Classification from MRI Data (DINOv2, BYOL & MoCo)	73.18%
Hybrid Approach for Autism Spectrum Disorder (ASD) Classification Using DINOv2, MoCo, and SimCLR Combined with EfficientNetB0 and ResNet50	Training: 93.5%, Test: 91.5%
ASD Classification Using DINOv2 and BYOL Models	93%
Leveraging DINOv2 and MoCo Models with ResNet50 for ASD Detection Using MRI Data	93%
Hybrid Model using DINOv2 and SimCLR for Autism Spectrum Disorder (ASD) Detection	98.01%

5. Conclusion

This study underscores the effectiveness of hybrid self-supervised learning frameworks in enhancing the detection of Autism Spectrum Disorder through MRI-based analysis. By combining state-of-the-art SSL models such as DINOv2, MoCo, SimCLR, and BYOL with CNN backbones, the proposed approaches achieved robust performance, with classification accuracies reaching up to 98.01%. The results validate that SSL-based pipelines can overcome the challenges of limited annotated datasets, improve generalization, and capture discriminative features essential for neuroimaging tasks (Baydar et al (2024) [23]). Beyond ASD classification, the methodology demonstrates scalability to broader applications in medical diagnostics, offering a pathway toward reliable, automated, and efficient tools for healthcare. Future work may expand this framework by incorporating multi-modal imaging, larger and more diverse datasets, and interpretability methods to further strengthen clinical trust and applicability.

In addition, future research should explore the integration of federated learning and privacy-preserving SSL techniques to enable collaborative training across multiple hospitals without sharing sensitive patient data (Jose et al (2025) [24]). Incorporating longitudinal neuroimaging data and clinical metadata, such as genetic markers and behavioral assessments, could provide a more holistic diagnostic framework. Furthermore, investigating explainable AI (XAI) methods within SSL-driven pipelines will be critical to ensure transparency, foster clinical acceptance, and build trust in AI-assisted ASD diagnosis (del

Olmo et al (2024) [25]). These directions can significantly advance the development of next-generation diagnostic systems that are accurate, interpretable, and ethically aligned with real-world healthcare practices.

References

- [1] Caron, M., H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin. “Emerging Properties in Self-Supervised Vision Transformers.” *Proceedings of the IEEE/CVF International Conference on Computer Vision*, (2021): 9650–9660.
- [2] Chen, T., S. Kornblith, M. Norouzi, and G. Hinton. “A Simple Framework for Contrastive Learning of Visual Representations.” *Proceedings of the International Conference on Machine Learning*, PMLR, (2020): 1597–1607.
- [3] He, K., H. Fan, Y. Wu, S. Xie, and R. Girshick. “Momentum Contrast for Unsupervised Visual Representation Learning.” *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, (2020): 9729–9738.
- [4] Doron, M., T. Moutakanni, Z. S. Chen, N. Moshkov, M. Caron, H. Touvron, ... and J. C. Caicedo. “Unbiased Single-Cell Morphology with Self-Supervised Vision Transformers.” *bioRxiv* (2023).
- [5] Tian, Y., L. Fan, K. Chen, D. Katabi, D. Krishnan, and P. Isola. “Learning Vision from Models Rivals Learning Vision from Data.” *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, (2024): 15887–15898.
- [6] Baydar, M. *Self-Supervised Learning for Unsupervised Image Classification and Supervised Localization Tasks*. Doctoral dissertation, Middle East Technical University (Turkey), 2024.
- [7] Xu, C., J. Wu, F. Zhang, J. Freer, Z. Zhang, and Y. Cheng. “A Deep Image Classification Model Based on Prior Feature Knowledge Embedding and Application in Medical Diagnosis.” *Scientific Reports* 14, no. 1 (2024): 13244.
- [8] Gaur, M., K. Chaturvedi, D. K. Vishwakarma, S. Ramasamy, and M. Prasad. “Self-Supervised Ensembled Learning for Autism Spectrum Classification.” *Research in Autism Spectrum Disorders* 107 (2023): 102223.
- [9] Wu, L., J. Zhuang, and H. Chen. “Voco: A Simple-Yet-Effective Volume Contrastive Learning Framework for 3D Medical Image Analysis.” *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, (2024): 22873–22882.
- [10] Grill, J. B., F. Strub, F. Alltech, C. Tallec, P. Richemond, E. Buchatskaya, and M. Valko. “Bootstrap Your Own Latent: A New Approach to Self-Supervised Learning.” *Advances in Neural Information Processing Systems* 33 (2020): 21271–21284.
- [11] Heinsfeld, A. S., A. R. Franco, R. C. Craddock, A. Buchweitz, and F. Meneguzzi. “Identification of Autism Spectrum Disorder Using Deep Learning and the ABIDE Dataset.” *NeuroImage: Clinical* 17 (2018): 16–23.

- [12] Baydar, M. *Self-Supervised Learning for Unsupervised Image Classification and Supervised Localization Tasks*. Doctoral dissertation, Middle East Technical University (Turkey), 2024.
- [13] Rani, V., S. T. Nabi, M. Kumar, A. Mittal, and K. Kumar. "Self-Supervised Learning: A Succinct Review." *Archives of Computational Methods in Engineering* 30, no. 4 (2023): 2761–2775.
- [14] Du, Y., J. Onofrey, and N. C. Dvornik. "Multi-View and Multi-Scale Alignment for Contrastive Language-Image Pre-Training in Mammography." *arXiv preprint arXiv:2409.18119* (2024).
- [15] Jose, C., T. Moutakanni, D. Kang, F. Baldassarre, T. Darcet, H. Xu, and P. Bojanowski. "DINOv2 Meets Text: A Unified Framework for Image-and Pixel-Level Vision-Language Alignment." *arXiv preprint arXiv:2412.16334* (2024).
- [16] Charoenpanyakul, R., V. Kittichai, S. Eiamsamang, P. Sriwichai, N. Pinetsuksai, K. M. Naing, and S. Chuwongin. "Enhancing Mosquito Classification through Self-Supervised Learning." *Scientific Reports* 14, no. 1 (2024): 27123.
- [17] Campanella, G., S. Chen, R. Verma, J. Zeng, A. Stock, M. Croken, B. Veremis, et al. "A Clinical Benchmark of Public Self-Supervised Pathology Foundation Models." *arXiv preprint arXiv:2407.06508* (2024).
- [18] Li, S., Y. Cao, Y. Liu, Z. Yan, Y. Dai, P. S. Yu, and L. Sun. "A Comprehensive Survey of AI-Generated Content (AIGC): A History of Generative AI from GAN to ChatGPT." *arXiv preprint arXiv:2303.04226* (2023).
- [19] Srivastava, D., V. Singh, S. Li, and K. Kochersberger. "Unmanned Aerial System-Driven Data and Advanced Deep Learning Strategies for Elevating Weed Management in Agriculture." (2024).
- [20] Morningstar, W., A. Bijamov, C. Duvarney, L. Friedman, N. Kalibhat, L. Liu, and S. Prakash. "Augmentations vs Algorithms: What Works in Self-Supervised Learning." *arXiv preprint arXiv:2403.05726* (2024).
- [21] Eapen, Naived George, and Jossy George. "Exploring Self-Supervised Learning Architectures for Image Processing: Milestones and Challenges." In *2024 IEEE 4th International Conference on ICT in Business Industry & Government (ICTBIG)*, pp. 1-5. IEEE, 2024.
- [22] Noreen, Muniba, and Furqan Shaukat. "Lung Nodule-SSM: Self-Supervised Lung Nodule Detection and Classification in Thoracic CT Images." *arXiv preprint arXiv:2505.15120* (2025).
- [23] Baydar, Melih. "Self-Supervised Learning for Unsupervised Image Classification and Supervised Localization Tasks." PhD diss., Middle East Technical University (Turkey), 2024.

- [24] Jose, Cijo, Théo Moutakanni, Dahyun Kang, Federico Baldassarre, Timothée Darcet, Hu Xu, Daniel Li et al. "Dinov2 meets text: A unified framework for image-and pixel-level vision-language alignment." In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 24905-24916. 2025.
- [25] del Olmo, Juan Jesus Losada, Ángel Luis Perales Gómez, Pedro E. Lopez-de-Teruel, and Alberto Ruiz. "A few-shot learning methodology for improving safety in industrial scenarios through universal self-supervised visual features and dense optical flow." *Applied Soft Computing* 167 (2024): 112375.
