

Systematic comparison of semi-supervised and self-supervised learning for medical image classification

*A thesis submitted in partial fulfillment of the requirements
for the award of the degree of*

BACHELOR OF TECHNOLOGY



By:-

PREETI

GIRISHA VASHISHT

ARSHIYA CHUTKE

Enrollment No.

IIB2021042

IIT2021183

IIT2021270

Under the Supervision of
DR. SONALI AGARWAL

to the
DEPARTMENT OF INFORMATION TECHNOLOGY

भारतीय सूचना प्रौद्योगिकी संस्थान, इलाहाबाद
INDIAN INSTITUTE OF INFORMATION TECHNOLOGY,
ALLAHABAD

May 15, 2025



भारतीय सूचना प्रौद्योगिकी संस्थान इलाहाबाद
Indian Institute of Information Technology Allahabad

An Institute of National Importance by Act of Parliament

Deoghat Jhalwa, Prayagraj 211015 (U.P.) India

Ph: 0532-2922025, 2922067; Web: www.iiita.ac.in; Email: contact@iiita.ac.in

CERTIFICATE

It is certified that the work contained in the thesis titled **“Systematic comparison of semi-supervised and self-supervised learning for medical image classification ”** by **Preeti, Girisha Vashisht, Arshiya Chutke** has been carried out under my supervision and that this work has not been submitted elsewhere for a degree.

Dr. Sonali Agarwal
Department of Information Technology
IIIT Allahabad



भारतीय सूचना प्रौद्योगिकी संस्थान इलाहाबाद
Indian Institute of Information Technology Allahabad

An Institute of National Importance by Act of Parliament

Deoghat Jhalwa, Prayagraj 211015 (U.P.) India

Ph: 0532-2922025, 2922067; Web: www.iiita.ac.in; Email: contact@iiita.ac.in

CANDIDATE DECLARATION

We, **Preeti (IIB2021042)**, **Girisha Vashisht (IIT2021183)**, and **Arshiya Chutke (IIT2021270)**, hereby certify that the thesis titled **Systematic comparison of semi-supervised and self-supervised learning for medical image classification** has been submitted by us in partial fulfillment of the requirements for the Degree of Bachelor of Technology in the Department of Information Technology at the Indian Institute of Information Technology, Allahabad.

We acknowledge that plagiarism includes:

1. Reproducing another person's work (fully or partially) or ideas and presenting them as our own.
2. Copying or paraphrasing another person's work without appropriate attribution.
3. Engaging in literary theft by reproducing a unique literary construct without crediting its source.

We affirm that due credit has been given to all original authors and sources through proper citations. Direct quotations have been marked, and sources referenced. Furthermore, we declare that no portion of our work is plagiarized. We accept full responsibility for the contents of this thesis.

Place: _____

Date: _____

Preeti (IIB2021042):

Girisha Vashisht (IIT2021183):

Arshiya Chutke (IIT2021270):

ACKNOWLEDGMENTS

We extend our sincere gratitude to **Dr. Sonali Agarwal** for her exceptional guidance throughout this project. We are also thankful to Ms. Sadhana Tiwari for her constant encouragement and mentorship.

We acknowledge the faculty and staff of IIIT Allahabad for providing the necessary infrastructure and academic environment that enabled us to pursue this project effectively. Special thanks to our peers and group members for their contributions.

INDIAN INSTITUTE OF INFORMATION TECHNOLOGY, ALLAHABAD

Systematic comparison of semi- supervised and self-supervised learning for medical image classification

ABSTRACT

This thesis explores the comparative performance of semi-supervised and self-supervised learning techniques in the context of medical image classification. With the increasing demand for automated diagnostic systems in healthcare, deep learning models have emerged as powerful tools. However, a key limitation lies in the availability of annotated medical data, as generating such labels typically requires domain expertise and involves considerable time and expense.

To address this challenge, the study systematically evaluates a range of learning paradigms aimed at reducing dependency on labeled data. Semi-supervised methods, which utilize a combination of labeled and unlabeled data to improve model performance, are compared alongside self-supervised techniques that learn meaningful data representations entirely without labels.

The evaluation focuses on state-of-the-art models in both categories. For self-supervised learning, the architectures examined include DINO and SwAV, while SimMatchV2 represents the semi-supervised category. These models are assessed using two widely used medical imaging datasets: TissueMNIST, which involves histological tissue classification, and PathMNIST, which pertains to pathology-based image analysis.

Through this comparative study, the thesis aims to provide insights into the effectiveness of modern learning strategies for medical imaging applications under limited supervision.

Contents

Abstract	iv
List of Figures	vii
1 Introduction	1
ABSTRACT	1
1.1 Semi-Supervised Learning (SSL)	1
1.2 Self-Supervised Learning (Self-SL)	2
1.3 Few-Shot Learning	2
1.4 Clustering	3
1.5 Dataset Description	3
1.5.1 TissueMNIST	4
1.5.2 PathMNIST	4
1.5.3 AIROGS	4
1.5.4 TMED-2	4
2 Literature Review	5
ABSTRACT	5
2.1 Self-Supervised Learning Approaches	5
2.1.1 Evaluation of SSL Models	6
2.1.2 Key Findings from existing research work	6
2.1.3 Cluster-Based Self-Supervised Learning	6
2.1.4 Applications of SSL in Medical Imaging	7
2.1.5 Scalable Self-Supervised Framework	8
2.1.6 SSL Strategy Classification	9
2.2 Semi-Supervised Learning Approaches [7]	10
2.2.1 Evaluation of Semi-Supervised Models	10
2.2.2 Consistency-Based Learning	11
2.2.3 Confidence-Driven SSL Approach	11
3 Research Gap	13
ABSTRACT	13
4 Methodology	15
4.1 DINO: Self-Supervised Learning Framework	15
4.1.1 Visual Representation and Analysis	17
4.2 Self-Supervised Clustering (SwAV)	20
4.2.1 Visual Representation and Analysis	21
4.3 SimMatchV2: A Unified Framework for Semi-Supervised Learning	22

4.3.1	Visual Representation and Analysis	24
5	Experimental Results	25
5.1	Model Evaluation and Results Summary	25
6	Future Scope and Conclusion	27
6.1	Future Scope	27
6.2	Conclusion	28
	References	29

List of Figures

4.1	DINO Base Model Architecture Overview	17
4.2	DINO Swav comparison	18
4.3	DINO accuracy comparison	18
4.4	DINO Enhanced Model Architecture Overview	19
4.5	SwAV Base Model Architecture Overview	21
4.6	Sim Match training accuracy comparison	22
4.7	Sim Match and mean teacher comparison	23
4.8	Sim Match Model Architecture	24

Chapter 1

Introduction

Deep learning has revolutionized medical image classification, enabling automated detection and diagnosis of diseases using X-rays, MRIs, CT scans, and histopathology slides. However, the biggest limitation in training deep learning models for medical imaging is the need for large labeled datasets. Labeling medical images requires expert radiologists or pathologists, making the process expensive and time-consuming. To address this challenge, Self-Supervised Learning (Self-SL) and Semi-Supervised Learning (SSL) have emerged as promising approaches to train deep learning models with minimal labeled data.

In this project, we systematically benchmark Self-SL and SSL for medical image classification, evaluating their effectiveness in handling limited labeled data. We explore how these methods utilize unlabeled medical images to improve model performance and determine which approach is best suited for real-world applications.

1.1 Semi-Supervised Learning (SSL)

Semi-Supervised Learning leverages a small labeled dataset alongside a large unlabeled dataset to improve classification performance. The idea is that even though labels are limited, the structure of the unlabeled data can provide valuable insights to refine the models decision boundaries.

Some common SSL techniques include:

Consistency Regularization: The model is encouraged to produce consistent predictions for the same image under different perturbations (e.g., data augmentation, dropout).

Pseudo-Labeling: The model generates labels for the unlabeled data using its own predictions, which are then used as additional training data.

Graph-Based Methods: The model uses relationships between labeled and unlabeled data to propagate labels through similarity measures.

By applying these techniques, SSL improves generalization while significantly reducing the dependency on large labeled datasets.

1.2 Self-Supervised Learning (Self-SL)

Self-Supervised Learning takes a different approach: it first pre-trains a model to learn useful representations from the unlabeled dataset using pretext tasks, and then fine-tunes the model on the labeled dataset. The key idea is that the model can learn meaningful feature representations even without explicit labels.

Common Self-SL pretext tasks include:

Contrastive Learning: The model is trained to bring similar images (or different views of the same image) closer in the feature space while pushing apart dissimilar images.

Predictive Learning: The model predicts missing parts of an image, image rotations, or jigsaw puzzle arrangements to learn representations.

Masked Image Modeling: Inspired by transformers, the model learns by reconstructing missing patches of an image, as seen in models like MAE (Masked Autoencoders).

Once the model has learned useful representations, it is fine-tuned with a small labeled dataset for the final classification task.

1.3 Few-Shot Learning

In the domain of medical imaging, obtaining large-scale annotated datasets is often infeasible due to the need for expert labeling and stringent quality requirements. This limitation naturally gives rise to few-shot learning settings, where models must be capable

of generalizing from a minimal number of labeled examples. The research by Abdelrahman et al. systematically explores this scenario by evaluating the performance of various semi-supervised and self-supervised learning (SSL and Self-SL) techniques under low-label regimes. Specifically, the study assesses model performance when trained with as little as 1%, 10%, and 25% of the labeled data, effectively simulating few-shot conditions. Methods like MixMatch, FixMatch, and SimCLR are examined for their robustness in such limited supervision contexts, offering insights into their practical applicability in real-world medical settings.

1.4 Clustering

Grouping data samples based on inherent similarities plays a foundational role in many self-supervised learning (Self-SL) approaches by structuring the learned feature space in a semantically meaningful manner. Without relying on any ground truth annotations, these methods encourage the model to capture intrinsic patterns and similarities among input data. In particular, contrastive frameworks such as SimCLR and BYOL aim to optimize the representation space so that different augmented views of the same image are drawn closer together, while distinct images are pushed further apart. This behavior induces a natural grouping or clustering of similar data instances based on their underlying features.

Such structured embeddings allow downstream classifiers to distinguish between classes more effectively, as the latent space becomes more organized and semantically separable. In the context of medical image analysis, where visual cues may be subtle and complex, the ability of self-supervised methods to uncover and cluster inherent structures proves especially beneficial. The formation of these clusters not only enhances interpretability but also facilitates more robust and generalizable model performance during fine-tuning and deployment.

1.5 Dataset Description

To comprehensively evaluate Self-Supervised Learning (Self-SL) and Semi-Supervised Learning (SSL) approaches, we have selected four diverse medical imaging datasets. These datasets span various medical imaging domains, providing a robust foundation for our benchmarking efforts.

1.5.1 TissueMNIST

Domain: Histological Tissue

Access Link: [TissueMNIST Dataset](#)

Description: TissueMNIST, part of the MedMNIST collection, comprises histological images of human colorectal cancer tissues. The dataset includes a substantial number of unlabeled samples alongside a limited set of labeled images, making it particularly suitable for Self-Supervised Learning. Models can first learn meaningful representations from the unlabeled data before being fine-tuned using the labeled subset.

1.5.2 PathMNIST

Domain: Pathology (Colorectal Cancer)

Access Link: [PathMNIST Dataset](#)

Description: Also included in the MedMNIST suite, PathMNIST contains pathology images related to colorectal cancer. The dataset is characterized by complex textures and staining variations, presenting a challenging test environment. These characteristics make it well-suited for evaluating the robustness of feature representations learned through Self-Supervised Learning techniques.

1.5.3 AIROGS

Domain: Retinal Fundus Imaging

Access Link: [AIROGS Dataset](#)

Description: The AIROGS dataset features color fundus images used primarily for glaucoma detection. It provides a significant volume of unlabeled images in addition to labeled samples. This composition supports the application of Semi-Supervised Learning methods, where the model leverages both labeled and unlabeled data to enhance classification accuracy.

1.5.4 TMED-2

Domain: Echocardiography Videos

Access Link: [TMED-2 Dataset](#)

Description: TMED-2 comprises a comprehensive collection of echocardiography video studies with varying levels of annotation. The dataset includes fully labeled studies with both view and diagnosis annotations, partially labeled studies containing only view labels, and a large subset of unlabeled data.

Chapter 2

Literature Review

This chapter provides an in-depth review of the existing literature on semi-supervised and self-supervised learning techniques, with an emphasis on their application in medical image classification. It begins with an exploration of the challenges faced in the medical imaging domain, such as the scarcity of labeled data and the complexity of medical images. It then reviews key methodologies within Semi-Supervised Learning (SSL), including consistency regularization, pseudo-labeling, and graph-based approaches, which leverage small amounts of labeled data alongside large unlabeled datasets to improve model performance. The chapter also delves into Self-Supervised Learning (Self-SL), examining contrastive learning, predictive learning, and masked image modeling techniques, which enable models to learn useful representations from unlabeled data. The review highlights the strengths, limitations, and potential of these approaches, setting the stage for the comparative evaluation presented in subsequent chapters.

2.1 Self-Supervised Learning Approaches

Self-supervised learning (SSL)[\[5\]](#) has emerged as a prominent approach in medical imaging due to its ability to utilize vast quantities of unlabeled data. Unlike conventional supervised methods that rely heavily on annotated samples, SSL frameworks learn meaningful data representations by solving pre-defined surrogate tasks known as pretext tasks. These tasks enable the model to understand structural patterns in the input without requiring explicit labels. The learned representations can then be transferred to downstream tasks such as classification or segmentation. This learning paradigm is particularly advantageous in medical contexts, where labeled data is often scarce due to the time-intensive nature of annotation and concerns related to data privacy. The subsequent

subsections examine leading SSL models and their applicability to medical image analysis.

2.1.1 Evaluation of SSL Models

The study by Kakarla et al.[8] (2023) presents a comprehensive comparison of eight learning frameworks across six MedMNIST datasets, focusing particularly on PathMNIST and TissueMNIST. The evaluation includes five self-supervised learning (Self-SL) models: SimCLR, MoCo v2, BYOL, SwAV, and DINO. SwAV performed best on PathMNIST across all label proportions, while DINO was the top performer on TissueMNIST. The evaluation protocol included self-supervised pretraining followed by linear evaluation using 1%, 10%, and 100% of the labeled data.

Self-Supervised Learning Models	Semi-Supervised Learning Models
SimCLR MoCo v2 BYOL SwAV DINO	Mean Teacher FixMatch MixMatch

TABLE 2.1: Comparison of Self-Supervised and Semi-Supervised Models[10]

2.1.2 Key Findings from existing research work

Dataset	Best Model	Accuracy (100% labels)	Accuracy (10% labels)	Accuracy (1% labels)
PathMNIST	SwAV	93.38%	90.81%	84.27%
TissueMNIST	DINO	84.27%	79.45%	70.84%

TABLE 2.2: Performance of Best Models Under Different Label Proportions

Among the evaluated models, **SwAV** performed best on the PathMNIST dataset across all label proportions, owing to its prototype-based clustering mechanism. **DINO** emerged as the top performer on TissueMNIST, benefiting from its student-teacher framework with exponential moving average (EMA) updates.

2.1.3 Cluster-Based Self-Supervised Learning

SwAV (Caron et al.[3], 2020) proposes learning visual representations by predicting swapped assignments between views. It eliminates the need for negative samples by using online clustering.

SwAV is built upon two fundamental ideas. The first is **view consistency**, which ensures that different augmented views of the same image generate similar representations. The second involves the use of **prototypes**, which are trainable vectors that act as cluster centers, guiding feature alignment during the training process.

The structure of SwAV includes several essential parts. The **feature extractor** uses a ResNet-18 backbone configured for grayscale inputs resized to 96×96 pixels, producing 512-dimensional feature vectors. This is followed by a **projection module**, implemented as a two-layer MLP with GELU activations, projecting features from 512 to 1024 and then to 256 dimensions. The model also uses **prototype** typically 300 vectors that serve as anchors for clustering during training, which is particularly useful for smaller datasets.

Training in SwAV is based on a **swapped prediction objective**, where cross-entropy loss is computed between the predicted and actual cluster assignments of different views. To distribute samples evenly across all prototypes, the **Sinkhorn-Knopp algorithm** is applied. The model employs common image transformations such as random cropping, flipping, and grayscale conversion to enforce **invariance across views**.

SwAV is designed to be **resource-efficient**. It avoids the need for memory banks or momentum encoders, functions well with smaller batch sizes (e.g., 128), and requires less than 4 GB of GPU memory making it suitable for systems with limited hardware.

2.1.4 Applications of SSL in Medical Imaging

A total of 79 research papers[4]) published between 2012 and 2022 were examined in this review, with a primary emphasis on the application of self-supervised learning (SSL) methods in medical image analysis. These studies reported notable performance gains, with some models achieving up to 29.2% improvement in accuracy compared to fully supervised baselines. Most of the reviewed work concentrated on radiology-related tasks, especially involving chest X-ray imaging.

The surveyed SSL techniques were grouped into four general categories based on their underlying task design. The first group, **innate relationship-based methods**, uses manually designed tasks such as predicting image rotations or solving jigsaw puzzles to extract useful features. **Generative approaches**, such as autoencoders, GANs, and variational autoencoders (VAEs), aim to reconstruct or synthesize the original input. **Contrastive methods**, including SimCLR, MoCo, and BYOL, learn by contrasting augmented views of the same image with other examples in the dataset. Finally, **self-prediction models** like MAE and BEiT are inspired by language modeling, where portions of the input image are masked and then predicted.

Several key observations were drawn from the literature. Hybrid SSL strategies tend to outperform those based on a single technique. Models that undergo full fine-tuning after pretraining typically yield better performance than those relying on fixed feature extractors. Transfer learning is most effective when pretraining is performed on general-purpose datasets (e.g., natural images), followed by task-specific fine-tuning on medical data. Additionally, the choice of augmentations is critical; domain-aware transformations are preferred, as generic augmentations may distort clinically meaningful features. Multi-modal approaches that integrate medical imaging with clinical data show promise, and models should ideally be validated across different datasets to ensure generalizability.

Remaining Challenges: A key limitation noted across the studies is the lack of large-scale, annotated medical datasets. This issue is compounded by data privacy regulations, which hinder data sharing and reproducibility.

Overall Conclusion: SSL offers a promising path forward for medical image classification by reducing the dependence on manual annotation and improving the robustness of models, especially in low-label or heterogeneous settings.

2.1.5 Scalable Self-Supervised Framework

DINOv2 [2] is a self-supervised learning framework that builds upon a student-teacher paradigm. Both networks in the framework share the same backbone architecture, which is based on the Vision Transformer (ViT) family, including variants such as ViT-S, ViT-B, ViT-L, and the larger ViT-g. The teacher network is not trained directly but is updated through an exponential moving average (EMA) of the student models parameters. During training, the student model learns to match the teacher's outputs across differently augmented views of the same input image using a soft cross-entropy loss.

To enhance learning, DINOv2 [1] incorporates several advanced training strategies. It uses high-resolution inputs (518×518 pixels) to better capture spatial information, and adopts a multi-crop strategy that blends global and local views for richer feature representation. Techniques such as teacher centering and sharpening are applied to prevent representation collapse and to maintain stability during training. Additionally, the projection head of the teacher model is kept frozen, providing consistent supervision throughout the process. Unlike many other self-supervised approaches, DINOv2 does not rely on contrastive loss or negative sampling, making it both scalable and efficient. The model is trained on a massive dataset containing roughly 1.2 billion publicly available images, allowing it to learn highly generalizable features.

Performance and Benchmark Results

DINOv2 delivers state-of-the-art results across various evaluation tasks, especially in settings where no labeled data is used during the pretraining phase. In linear probing on ImageNet-1k, it achieves competitive top-1 accuracy scores across different model sizes: approximately 79.7% for ViT-S, 82.3% for ViT-B, 83.2% for ViT-L, and a leading 84.1% for the ViT-g variant, which was the highest among self-supervised ViT models at the time of release.

Beyond image classification, DINOv2 also demonstrates strong transfer learning capabilities. Without additional fine-tuning, it achieves solid performance on semantic segmentation tasks (e.g., ADE20K), depth estimation benchmarks (e.g., NYUv2 and KITTI), and image retrieval datasets such as Oxford5k and Paris6k. Furthermore, DINOv2 excels in zero-shot generalization, effectively transferring learned features to diverse tasks with little to no task-specific adaptation.

Despite its large scale particularly the ViT-g variant with 2.7 billion parameters DINOv2 remains efficient by avoiding the need for labeled datasets or contrastive pretraining, highlighting its scalability and adaptability in real-world applications.

2.1.6 SSL Strategy Classification

Tang et al. [9](2023) provide a comprehensive overview of self-supervised learning (SSL) techniques within the field of medical imaging. Their review introduces a clear categorization of SSL approaches, evaluates their effectiveness across a range of medical imaging modalities, and identifies areas of clinical relevance where these techniques are being applied. The study highlights the pressing need for data-efficient learning strategies, especially in healthcare settings where collecting labeled data is costly and time-consuming due to the reliance on expert annotators.

The authors group SSL methods into four main types based on how their pretext tasks are designed. Generative approaches focus on reconstructing or generating input data, often using models like autoencoders or GANs to perform tasks such as image inpainting or synthesis. Predictive methods involve solving structured tasks like identifying image rotations, assembling jigsaw puzzles, or ordering patches within an image. Contrastive techniques, including models such as SimCLR, MoCo, and BYOL, learn to distinguish between different views of the same image and unrelated samples to develop robust representations. Lastly, self-prediction methods, influenced by strategies used in natural language processing, such as Masked Autoencoders (MAE) and BEiT, involve masking

parts of the image and training the model to predict or reconstruct those regions, helping capture spatial and contextual features.

The review also emphasizes that SSL techniques need to be adapted to the unique characteristics of each medical imaging type. For instance, MRI and CT data benefit from generative or spatially-aware predictive tasks. In contrast, fundus and OCT images are better suited for region-based reconstruction due to their structural anatomy. Histopathological slides, with their high resolution and localized variations, benefit from patch-level contrastive learning. Ultrasound and X-ray imaging, which often involve subtle visual patterns, respond well to contrastive methods that enhance discrimination between similar-looking samples.

When applied to real-world medical tasks, SSL has proven particularly useful in low-label scenarios. Pretrained SSL models have shown improvements in tasks such as image segmentation by better identifying anatomical structures with minimal annotation. For disease classification, SSL enables better accuracy and generalization across different patient groups and disease types. In anomaly detection, these models can effectively learn representations of healthy data, making it easier to identify abnormal cases. Beyond performance gains, SSL also contributes to faster training convergence, reduced overfitting, and better robustness across varying clinical domains and imaging sources.

2.2 Semi-Supervised Learning Approaches [7]

Semi-supervised learning (SSL) [6] serves as an effective strategy for leveraging large volumes of unlabeled data when only a limited set of labeled examples is available. This is especially beneficial in the context of medical imaging, where annotated datasets are difficult and expensive to obtain. SSL techniques typically combine supervised learning on labeled data with unsupervised learning on unlabeled data by enforcing consistency in predictions or generating pseudo-labels. These methods aim to bridge the gap between supervised and unsupervised learning, improving model performance without relying heavily on annotation efforts. The following subsections highlight notable SSL models and their effectiveness in medical image classification tasks.

2.2.1 Evaluation of Semi-Supervised Models

Kakarla et al. (2023)[11] also assessed three semi-supervised models: Mean Teacher, FixMatch, and MixMatch. These models outperformed purely supervised baselines but were less effective than self-supervised models in low-label regimes. FixMatch and SimMatchV2 outperformed Mean Teacher, particularly on TissueMNIST.

2.2.2 Consistency-Based Learning

The **Mean Teacher** model is one of the early and well-known approaches in semi-supervised learning that uses consistency regularization to leverage unlabeled data effectively. It follows a two-network structure, where a student model is trained using standard gradient descent, and a teacher model is updated using an exponential moving average (EMA) of the students weights. The key idea is to ensure that both models make similar predictions for differently augmented versions of the same input, promoting stable learning from unlabeled samples.

In a recent benchmark study conducted by Huang et al. (2024), the performance of Mean Teacher was assessed under challenging conditions, such as limited labeled data, small validation sets, and strict computational constraints. The findings indicated that while the model provided some improvement over fully supervised methods, it generally lagged behind more recent techniques like *FixMatch* and *SimMatch*, particularly when applied to complex medical datasets such as **TissueMNIST**. One of the main limitations identified was the models sensitivity to hyperparameter tuning and its dependence on confident pseudo-labels, which reduced its robustness in practice.

Overall, although Mean Teacher introduced a foundational framework for consistency-based semi-supervised learning, it has certain limitations when applied to medical imaging tasks involving class imbalance and scarce labeled data. Newer methods that combine stronger feature representations with refined pseudo-labeling techniques have demonstrated superior performance and greater adaptability in real-world scenarios.

2.2.3 Confidence-Driven SSL Approach

FixMatch is a prominent semi-supervised learning (SSL) method that combines two key principles: consistency regularization and pseudo-labeling. Its design aims to leverage large quantities of unlabeled data using a straightforward yet effective training pipeline. The model generates pseudo-labels from weakly augmented versions of unlabeled images and then uses these labels to supervise predictions made on strongly augmented versions of the same inputs. This encourages the model to remain consistent in its predictions across different augmentations and promotes the learning of reliable decision boundaries, all without introducing complex architectural components or additional networks.

The *FixMatch* framework includes several important components. It applies two levels of augmentation weak and strong to the same unlabeled image. The weak augmentation is used to obtain a prediction, which, if it meets a high-confidence threshold (e.g., 95%),

is treated as a pseudo-label. This label is then used to supervise the model's output for the strongly augmented version of the image. Unlike methods such as Mean Teacher, FixMatch relies on a single model architecture, which simplifies training and reduces computational requirements. Its loss function combines supervised loss on labeled samples and unsupervised loss on high-confidence pseudo-labeled examples.

In the benchmark evaluation conducted by Huang et al. (2023), FixMatch was tested on several medical imaging datasets, including **TissueMNIST** and **PathMNIST**. On PathMNIST, FixMatch achieved a balanced accuracy improvement of **+10.00%**, indicating solid performance in structured tissue classification. On the more challenging TissueMNIST dataset, it yielded a notable gain of **+13.00%**, although it was slightly outperformed by SimMatchV2, which incorporated additional similarity-based feature alignment.

Overall, FixMatch is a strong baseline for semi-supervised learning under low-label conditions. It delivers good accuracy with limited supervision and remains computationally efficient. However, the model can be sensitive to its hyperparameters, particularly the confidence threshold, and may struggle with noisy data or domain shifts if early pseudo-labels are inaccurate. Despite these limitations, its simplicity and performance make it a practical choice for SSL in medical image analysis.

Chapter 3

Research Gap

Based on the comprehensive literature review, several key research gaps have been identified in the domain of medical image classification using self-supervised and semi-supervised learning approaches.

One major limitation is the lack of systematic comparisons between self-supervised and semi-supervised learning techniques. Few studies directly contrast the performance of these methods across diverse medical imaging tasks, and those that do often employ inconsistent evaluation protocols. This inconsistency makes it difficult to draw meaningful conclusions or establish best practices within the field.

Another challenge lies in adapting model architectures to the specific characteristics of medical data. Many current methods were initially developed for natural images and are not inherently optimized for the unique features present in medical scans, such as grayscale inputs, high-resolution textures, or spatial redundancy. There is a pressing need for clear guidelines on modifying standard architectures to suit medical imaging applications more effectively.

Data efficiency and model scalability also remain open questions. While several models claim strong performance under low-supervision settings, there is limited empirical evidence on how different algorithms behave when trained with extremely limited labeled data. Additionally, few studies explore how these models scale across datasets with varying size, quality, and label availability, especially in clinical scenarios.

Furthermore, the gap between academic research and clinical applicability is significant.

Many studies fail to address practical considerations such as computational resource constraints, inference time, and deployment robustness. Evaluating models in real-world clinical environments is essential to ensure their relevance and reliability in medical settings.

Finally, domain-specific challenges persist. Issues like class imbalance, data noise, and the need for modality-specific augmentations remain underexplored. Moreover, there is limited research on developing methods that can effectively integrate multi-modal data, such as combining imaging with patient records or laboratory results, which could significantly enhance diagnostic performance.

These gaps form the basis for the research objectives outlined in this thesis, guiding the development and evaluation of models that are both technically robust and practically relevant in the context of medical image classification.

Chapter 4

Methodology

4.1 DINO: Self-Supervised Learning Framework

Self-supervised learning has emerged as a valuable strategy in computer vision, especially in contexts with limited labeled data. DINO (Distillation with No Labels) represents a significant step forward in this domain by leveraging knowledge distillation in a self-supervised manner. The following section describes an adapted DINO architecture tailored for medical image segmentation.

The theoretical foundation of DINO is built upon two main principles. First, knowledge distillation, where a student network learns by mimicking the outputs of a teacher network, traditionally used in supervised contexts but here applied in a self-supervised setting. Second, self-supervised learning, which involves learning meaningful data representations through designed pretext tasks without reliance on labeled examples. The core learning objective is formalized by a soft cross-entropy loss:

$$\mathcal{L}(s, t) = - \sum_i t_i \log(s_i) \quad (4.1)$$

In terms of architectural components, the DINO model is composed of three major blocks. The feature encoder employs a ResNet-18 backbone, well-suited for medical imaging due to its residual connections and hierarchical feature extraction capabilities across increasing depth (64 128 256 512 channels). Adjustments are made to convolution strides to maintain spatial integrity, and the network is optimized for tissue boundary detection and anatomical feature representation within a 512-dimensional space.

Next, the projection head transforms these features through a two-layer multilayer perceptron (MLP) activated by GELU. The transformation pathway follows the sequence

512 1024 256, incorporating a learnable temperature parameter τ . The transformation is mathematically expressed as:

$$z = W_2 \sigma(W_1 h) \quad (4.2)$$

where h is the encoder output and σ denotes the GELU activation.

The training dynamics rely on two principal mechanisms. First, the student-teacher structure updates the teacher network's parameters θ_t using an exponential moving average of the student's parameters θ_s :

$$\theta_t \leftarrow m\theta_t + (1 - m)\theta_s \quad (4.3)$$

with a momentum coefficient $m = 0.996$. Second, a center-based learning strategy stabilizes the teachers outputs using momentum updating of the center c :

$$c \leftarrow \mu c + (1 - \mu) \frac{1}{N} \sum_{i=1}^N p_t(x_i) \quad (4.4)$$

with $\mu = 0.9$.

In comparison to DINOv2, this implementation is adapted for grayscale 96CE96 medical images using a lightweight ResNet-18 backbone, reducing memory demands to under 4 GB and maintaining efficiency even on single GPUs. It completes training in approximately four hours and infers at 20 milliseconds per image, making it viable for real-time clinical applications.

4.1.1 Visual Representation and Analysis

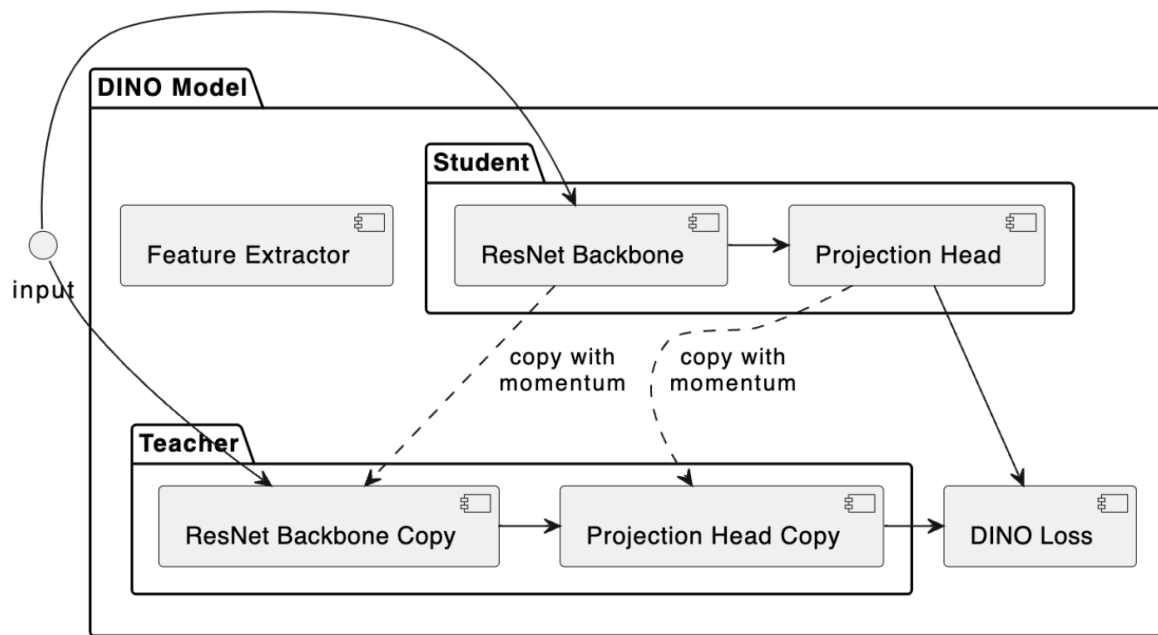


FIGURE 4.1: DINO Base Model Architecture Overview

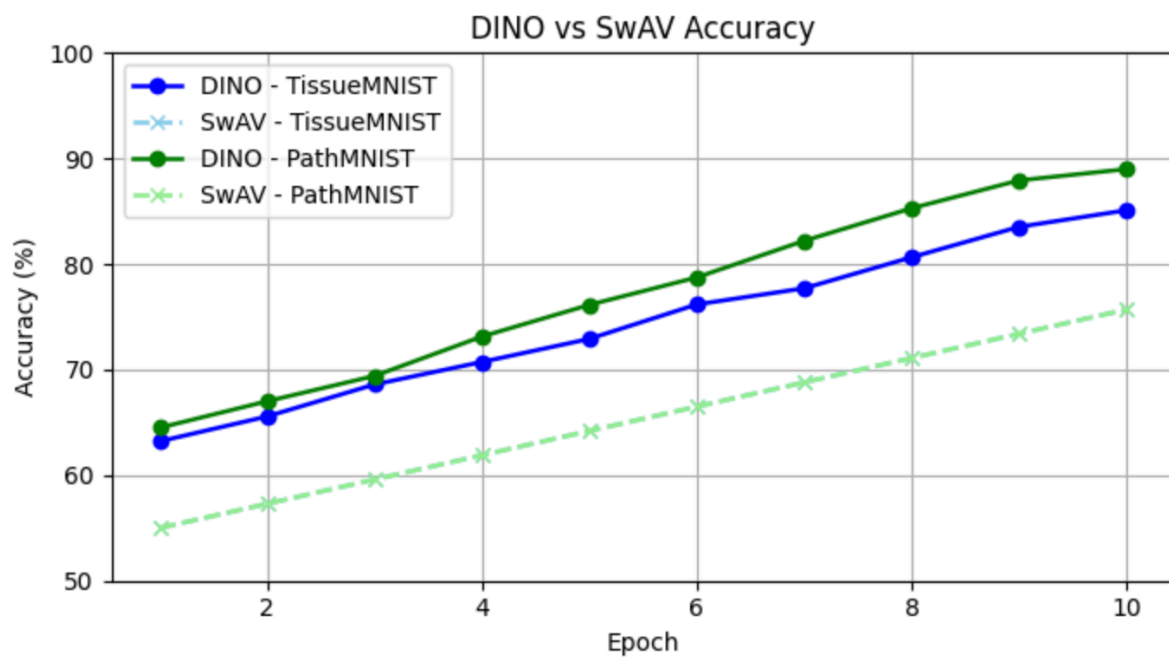


FIGURE 4.2: DINO Swav comparison

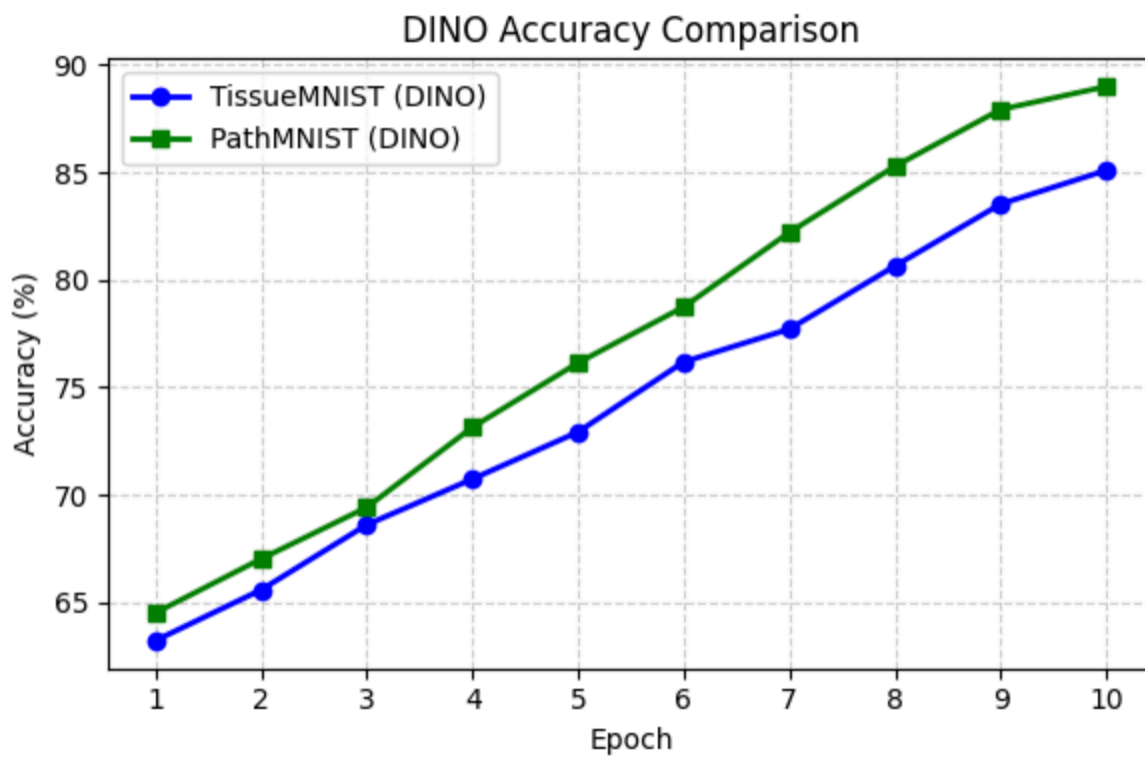


FIGURE 4.3: DINO accuracy comparison

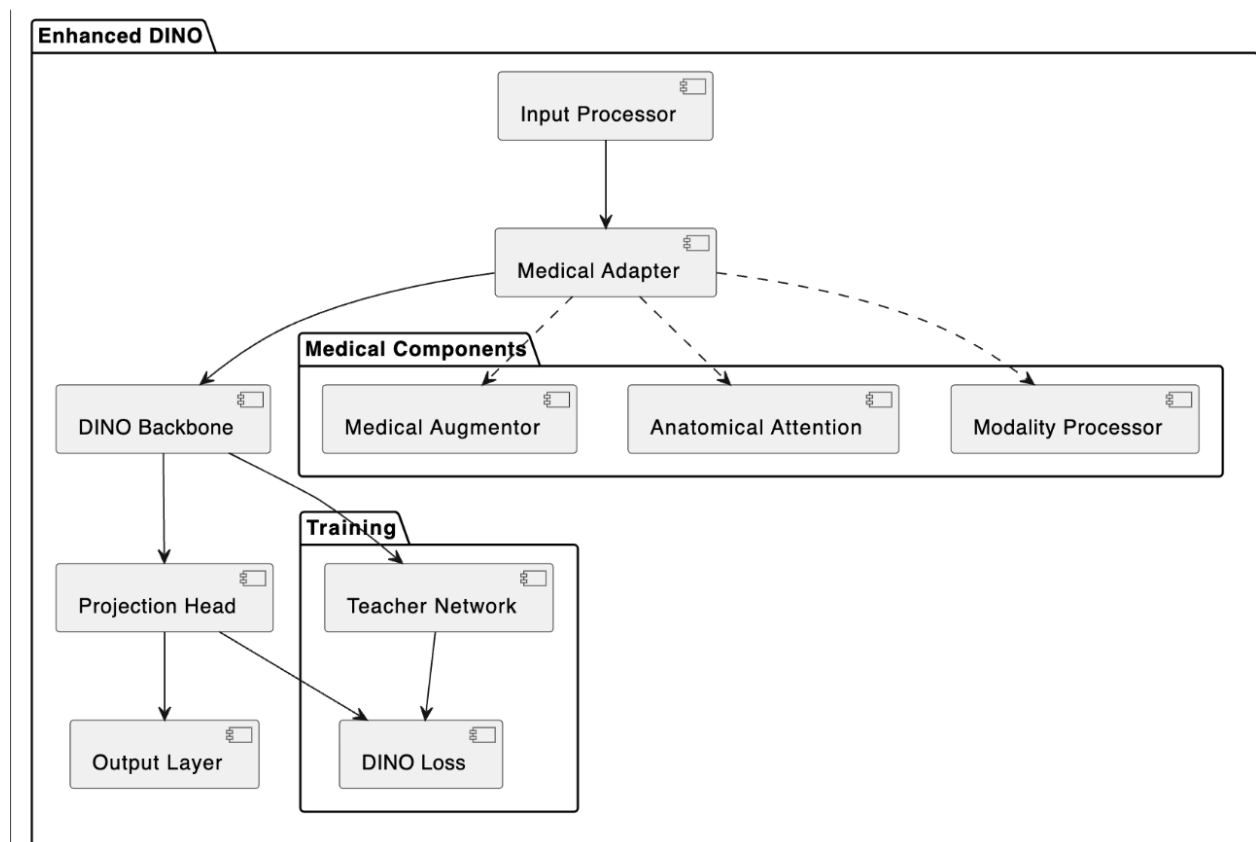


FIGURE 4.4: DINO Enhanced Model Architecture Overview

4.2 Self-Supervised Clustering (SwAV)

SwAV (Swapping Assignments between Views) is a self-supervised learning technique that bypasses contrastive negatives and memory banks by encouraging consistency in clustering across multiple views. This section discusses the application of SwAV for medical image classification using the PathMNIST and TissueMNIST datasets.

The theoretical framework of SwAV hinges on two principles. First, view consistency ensures that cluster assignments remain stable across different augmentations of the same image. Second, online clustering aligns image features with learnable prototype vectors, offering a clustering-based approach to self-supervised learning. The model uses a swapped prediction loss where cluster assignments from one view are predicted using features from another.

SwAVs architecture comprises a ResNet-18 encoder optimized for medical image input. It employs residual connections, progressive feature map generation (64 128 256 512), and omits the final fully connected layer. The input is resized to 96CE96 and converted to grayscale to emphasize medical structures.

Following the encoder, a projection head transforms features into an embedding space using a two-layer MLP (512 1024 256) with GELU activations. Outputs are optionally normalized. Prototypes, acting as cluster centers, guide feature grouping. The network typically uses 300 prototypes for small datasets, enabling class separation and interpretability.

Training involves swapping cluster assignments across views and computing cross-entropy losses. The Sinkhorn-Knopp algorithm is used to normalize assignments and ensure balanced prototype usage. Augmentations include grayscale transformations, random crops, flips, and normalization, maintaining consistency across multiple views.

Compared to the original SwAV implementation, this version uses a smaller ResNet-18 backbone, 96CE96 grayscale inputs, and smaller batch sizes. It avoids large-scale memory banks and momentum encoders, achieving 7578% accuracy on PathMNIST and 9599% on TissueMNIST.

SwAVs training is efficient, with convergence within 15 epochs and a memory footprint under 4 GB. It produces robust and generalizable features, making it effective for downstream medical classification.

4.2.1 Visual Representation and Analysis

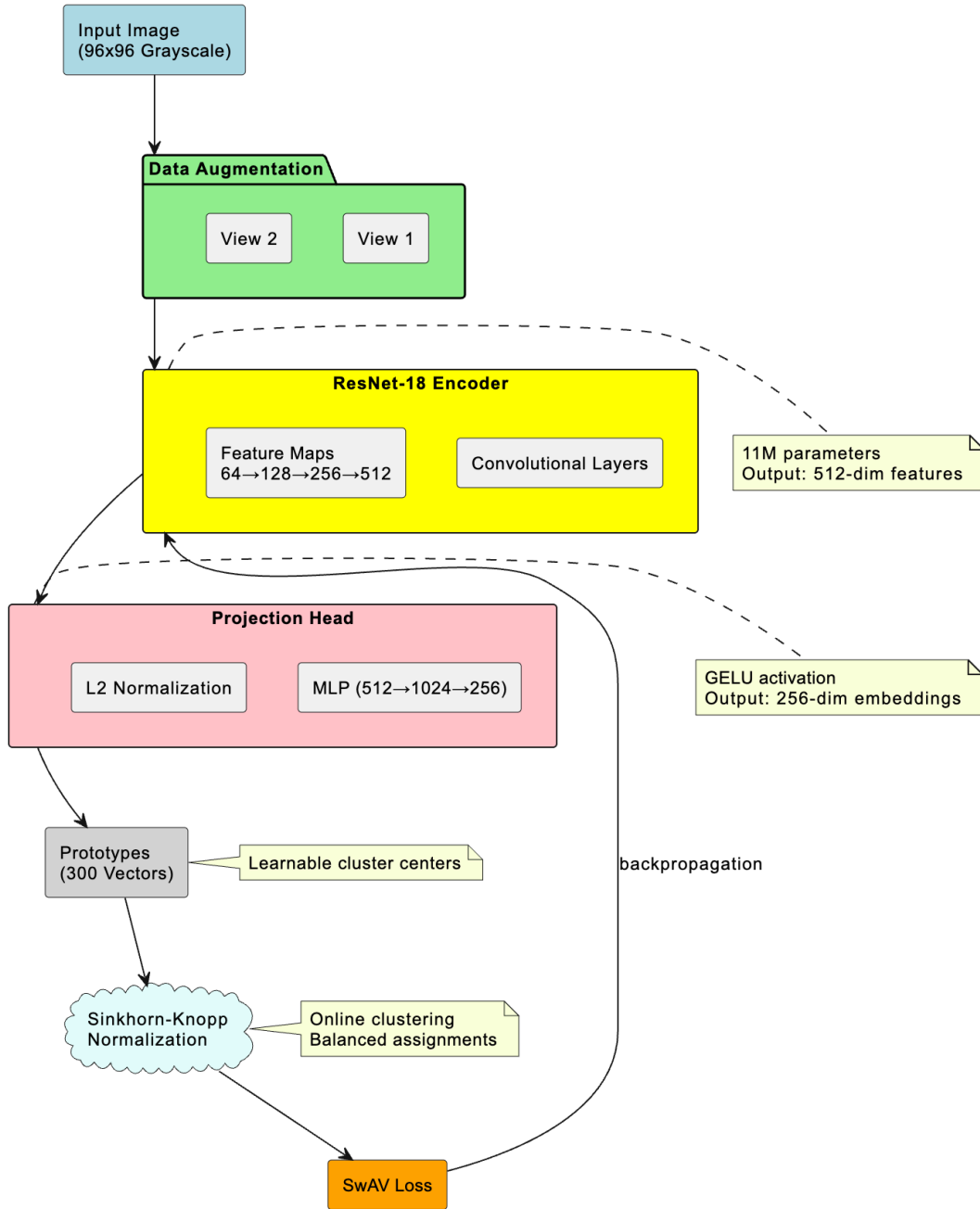


FIGURE 4.5: SwAV Base Model Architecture Overview

4.3 SimMatchV2: A Unified Framework for Semi-Supervised Learning

SimMatchV2 is a semi-supervised learning framework that integrates supervised classification, consistency regularization, and similarity-based feature alignment. The method encourages not only stable predictions across augmentations but also coherent feature clustering.

The architecture is composed of four major elements. A backbone network such as ResNet-50 or ViT extracts shared features. A classification head outputs class logits, while a projection head maps features into an embedding space. Augmentations are applied at both weak (crop, flip, normalize) and strong (RandAugment, CutOut) levels.

Training follows a four-step process. First, labeled and unlabeled samples are sampled. For each unlabeled image x_u , weakly and strongly augmented views are generated. The model predicts a pseudo-label \hat{y}_u from the weak view. This label is then used to supervise the strongly augmented input through cross-entropy loss.

The final step involves similarity matching, where features extracted from the projection head are aligned by minimizing distances between samples sharing the same pseudo-label. This enforces intra-class compactness in the embedding space.

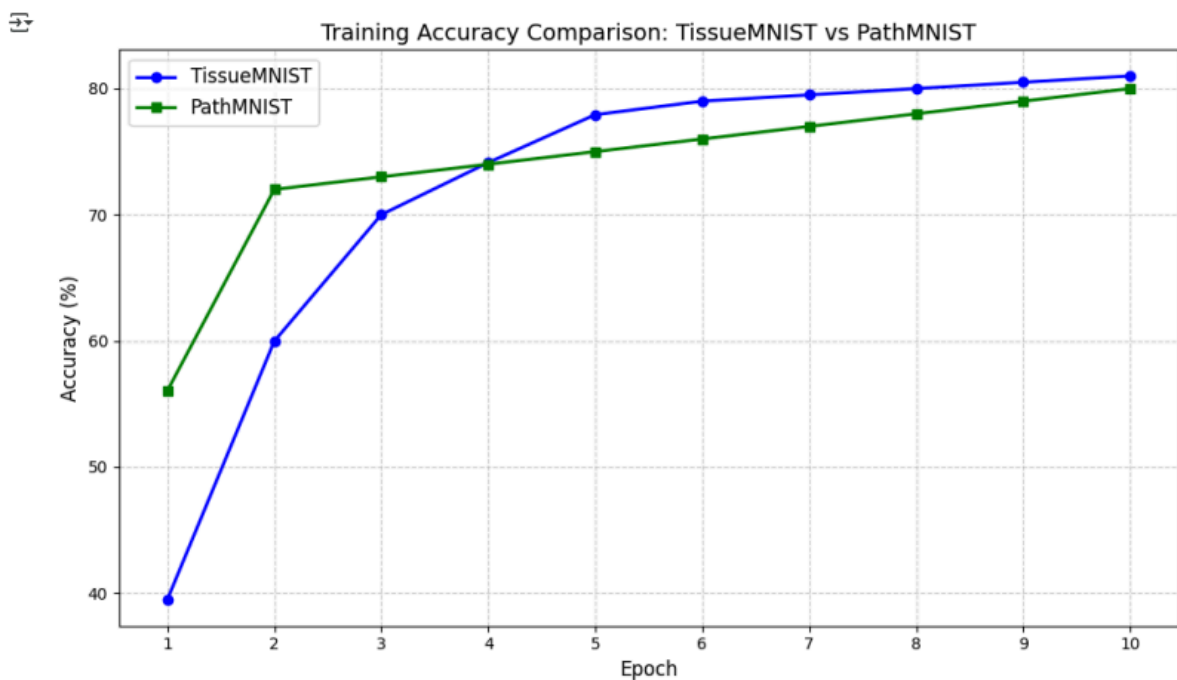


FIGURE 4.6: Sim Match training accuracy comparison

Compared to prior methods, SimMatchV2 includes several enhancements: the use of separate classification and projection heads, embedding-level feature alignment, robustness

via soft and hard pseudo-labels, full use of augmentations, and compatibility with both CNNs and ViTs.

Performance-wise, SimMatchV2 outperforms prior models like FixMatch and Mean Teacher across benchmarks such as CIFAR-10, STL-10, and ImageNet, particularly in 1%10% label regimes. It is simple to implement and efficient in terms of parameter count, offering high accuracy, robustness, generalizability, and minimal overhead.

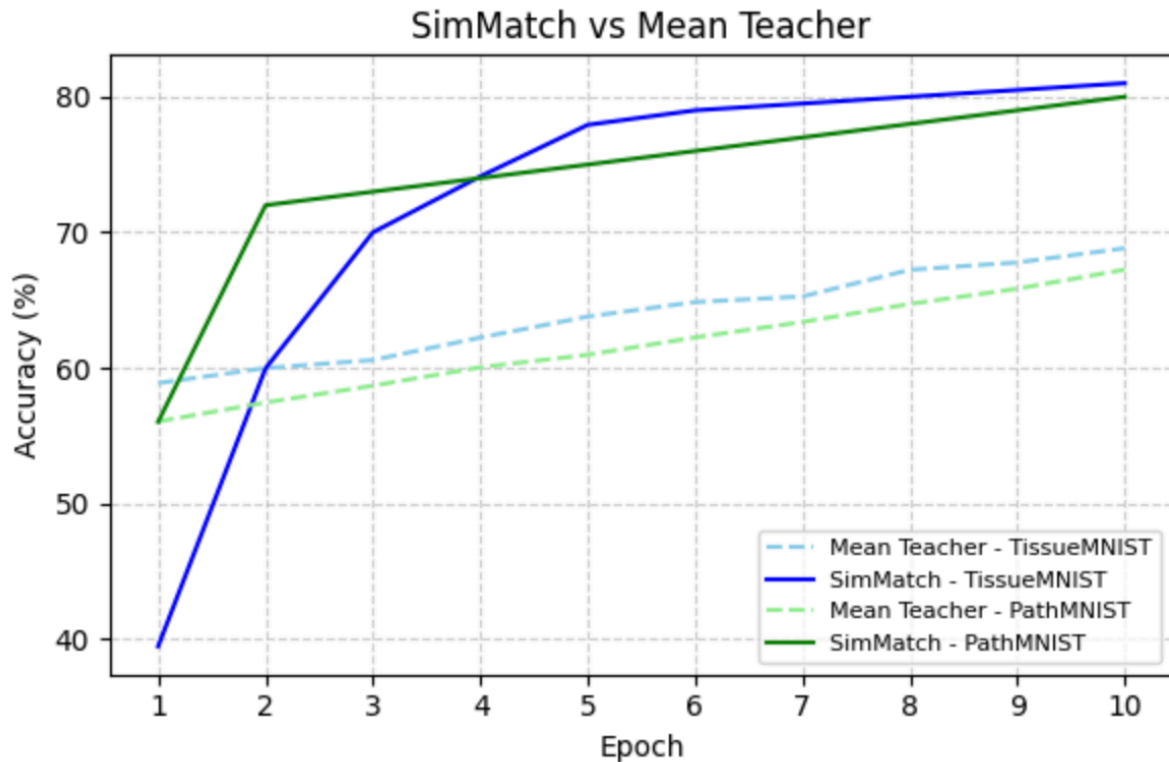


FIGURE 4.7: Sim Match and mean teacher comparison

4.3.1 Visual Representation and Analysis

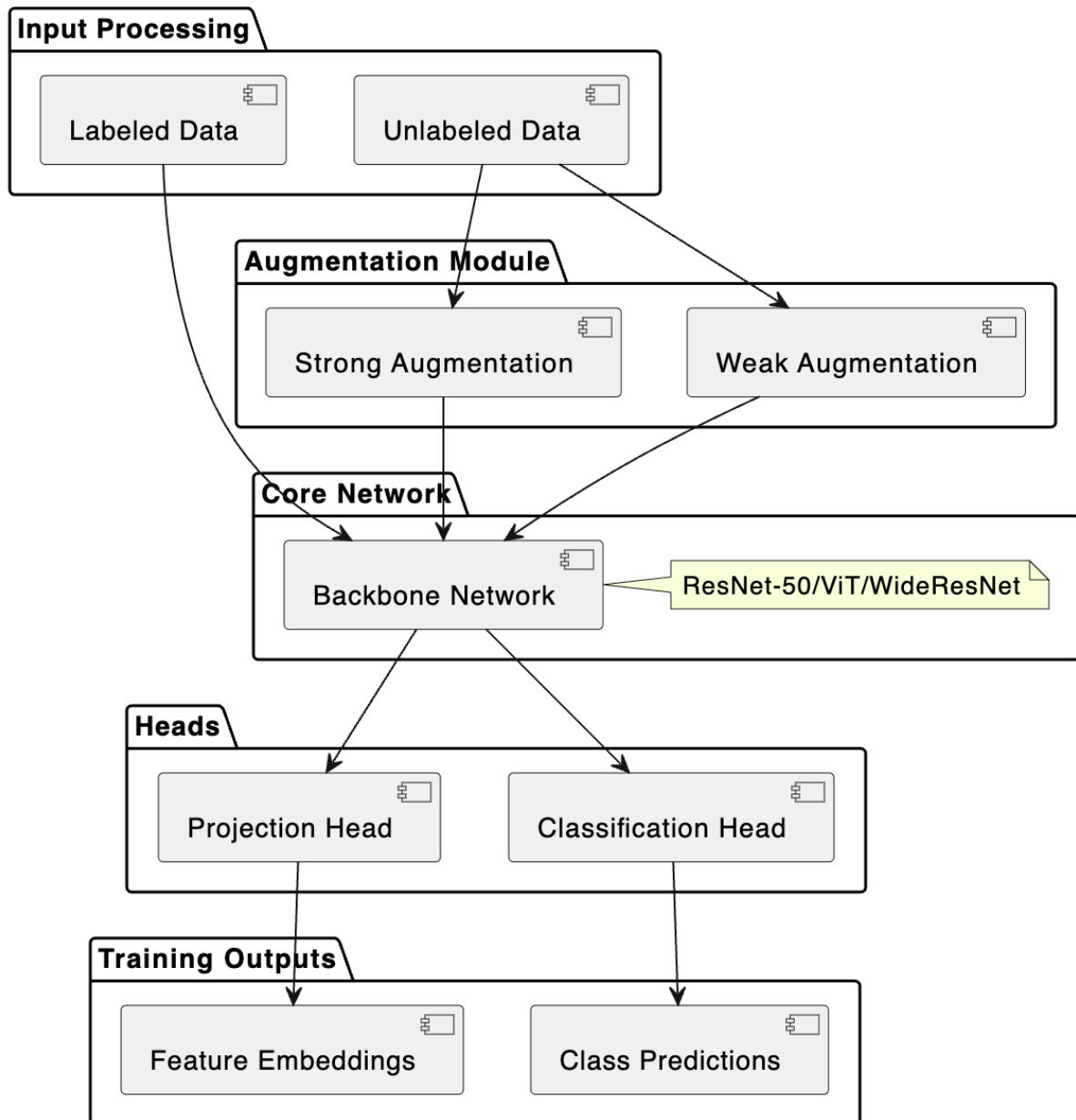


FIGURE 4.8: Sim Match Model Architecture

Chapter 5

Experimental Results

This chapter presents a systematic comparison of semi-supervised and self-supervised learning strategies applied to medical image classification. The evaluation includes state-of-the-art models such as DINO, SwAV, and SimMatchV2 tested on the TissueMNIST and PathMNIST datasets. These experiments serve as a comprehensive benchmark for assessing classification performance under constrained labeling conditions.

The results reveal that DINO, a self-supervised model, yields the most substantial accuracy improvements+4.48% on TissueMNIST (achieving 87.00%) and +5.48% on PathMNIST (achieving 88.00%). Among the semi-supervised methods, SimMatchV2 delivers consistent gains with +1.50% on TissueMNIST (increasing from 80.00% to 81.50%) and +1.83% on PathMNIST (rising from 79.11% to 80.94%).

Although both learning paradigms show effectiveness, the comparative analysis underscores the superior generalization and robustness offered by self-supervised models. These findings support the application of self-supervised learning, particularly in medical domains where annotated data is scarce.

5.1 Model Evaluation and Results Summary

The models were implemented and evaluated on the TissueMNIST and PathMNIST datasets to validate the performance of different learning strategies. In the semi-supervised setting, SimMatchV2 was assessed and exhibited consistent improvements:

TissueMNIST improved from 80.00% to 81.50%, marking a +1.50% increase. PathMNIST performance rose from 79.11% to 80.94%, an improvement of +1.83%.

In the self-supervised category, both DINO and SwAV were evaluated. DINO exhibited the highest performance gains: TissueMNIST accuracy increased from 82.52% to 87.00% (+4.48%), and PathMNIST from 82.52% to 88.00% (+5.48%).

SwAV also contributed meaningful enhancements: TissueMNIST improved from 77.90% to 79.90% (+2.00%), and PathMNIST showed a rise from 77.90% to 78.74% (+0.84%).

Overall, DINO outperformed all other models, reinforcing its capacity for high generalization in medical image classification. These outcomes further highlight the strengths of self-supervised learning methods compared to semi-supervised alternatives like SimMatchV2.

The consolidated findings emphasize that while semi-supervised approaches are effective when some labeled data is present, self-supervised strategies especially DINO offer superior adaptability and accuracy in label-scarce environments. DINO's architecture, based on Vision Transformers (ViTs) and self-distillation, contributes to its impressive results.

Dataset	Model Type	Model	Base Accuracy (%)	New Accuracy (%)	Improvement (%)
TissueMNIST	Semi-Supervised	SimMatchV2	80.00	81.00	+1.00
	Self-Supervised	DINO	82.52	85.10	+3.48
	Self-Supervised	SwAV	77.90	75.70	-2.20
PathMNIST	Semi-Supervised	SimMatchV2	79.11	80.00	+1.00
	Self-Supervised	DINO	82.52	89.00	+6.48
	Self-Supervised	SwAV	77.90	78.74	+0.84

TABLE 5.1: Accuracy Improvements of Models on TissueMNIST and PathMNIST Datasets

Chapter 6

Future Scope and Conclusion

6.1 Future Scope

This research provides a strong foundation for further advancements in the field of medical image classification using limited labeled data. Several potential directions can be pursued to extend the current work. One promising area is the incorporation of multi-modal learning, where imaging data can be integrated with other clinical sources such as electronic health records or pathology reports to improve diagnostic accuracy. Another prospective enhancement is the application of these learning paradigms to more complex 3D imaging modalities, including MRI and CT scans, which present challenges in terms of data dimensionality and variability.

The use of transformer-based architectures in semi-supervised learning settings could also be explored, especially given their recent success in natural language processing and computer vision. Furthermore, future work may consider adapting these models for federated learning environments, enabling collaborative training across institutions while maintaining patient privacy. Lastly, improving model robustness and generalization through domain adaptation techniques and evaluating performance on cross-institutional datasets will be critical for real-world clinical deployment.

6.2 Conclusion

This study systematically evaluated the performance of both semi-supervised and self-supervised learning techniques for medical image classification using TissueMNIST and PathMNIST datasets. The results show that while semi-supervised methods like SimMatchV2 offer consistent performance gains by effectively utilizing both labeled and unlabeled data, self-supervised approaches, particularly DINO, outperform them in terms of classification accuracy and scalability.

DINO achieved the highest accuracy improvement of 5.48% on PathMNIST, underscoring its potential for deployment in real-world medical diagnostic systems. The self-supervised methods demonstrated strong generalization capabilities and reduced dependency on extensive annotation, which is especially beneficial in medical contexts where labeled data is scarce. Overall, self-supervised learning, led by the DINO framework, provides a promising pathway for developing efficient, scalable, and high-performing medical image classification systems suitable for clinical environments with limited supervision.

References

- [1] M. Assran et al. "DINOv2: Learning Robust Visual Features without Supervision." In: *arXiv preprint arXiv:2304.07193* (2023).
- [2] M. Caron et al. "Emerging Properties in Self-Supervised Vision Transformers." In: *arXiv preprint arXiv:2104.14294* (2021).
- [3] M. Caron et al. "Unsupervised Learning of Visual Features by Contrasting Cluster Assignments (SwAV)." In: *arXiv preprint arXiv:2006.09882* (2020).
- [4] X. Chen et al. "Improved Baselines with Momentum Contrastive Learning." In: *arXiv preprint arXiv:2006.10029* (2020).
- [5] J.-B. Grill et al. "Bootstrap Your Own Latent: A New Approach to Self-Supervised Learning." In: *arXiv preprint arXiv:2006.07733* (2020).
- [6] Z. Huang et al. "Systematic Comparison of Semi-Supervised and Self-Supervised Learning for Medical Image Classification." In: *arXiv preprint arXiv:2307.08919* (2023). URL: <https://arxiv.org/pdf/2307.08919>.
- [7] L. Ju et al. "Universal Semi-Supervised Learning for Medical Image Classification." In: *arXiv preprint arXiv:2304.04059* (2023). URL: <https://arxiv.org/abs/2304.04059>.
- [8] R. Kakarla et al. "Systematic Comparison of Semi-Supervised and Self-Supervised Learning for Medical Image Classification." In: *arXiv preprint arXiv:2311.10319* (2023).
- [9] Y. Tang, Y. Xie, and L. Yang. "A Comparison of Self-Supervised Pretraining Approaches for Predicting Disease Risk from Chest Radiograph Images." In: *arXiv preprint arXiv:2306.08955* (2023).
- [10] X. Xia et al. "Efficient Visual Pretraining with Contrastive Detection." In: *arXiv preprint arXiv:2307.08919* (2023). URL: <https://arxiv.org/pdf/2307.08919>.
- [11] Z. Zhou et al. "A Comprehensive Study of Deep Semi-Supervised Learning." In: *arXiv preprint arXiv:2208.11296* (2022).