# Few-shot and Transfer Learning for Deep Learning in Medical Imaging: Improving Abnormality Detection with Data-Efficient and Generative Approaches

Krish Patel
*Department of Computer Science*
*University of Illinois Chicago*
kpate446@uic.edu

Rohit Nanjundareddy
*Department of Computer Science*
*University of Illinois Chicago*
rnanj@uic.edu

*Abstract*—Medical imaging plays a critical role in clinical diagnosis and treatment planning, yet the development of robust deep learning models is often constrained by limited annotated data, high annotation costs, and stringent privacy regulations. This report presents a comprehensive investigation of few-shot learning and transfer learning techniques for medical abnormality detection. We evaluate multiple pre-trained architectures including ResNet-18, Vision Transformers (ViT), and medical-specific models (DenseNet-121/CheXNet) on the CheXpert chest X-ray dataset. Our experiments demonstrate that few-shot learning with domain adaptation significantly outperforms zero-shot baselines, achieving up to 48.5% accuracy in 5-shot scenarios compared to 38.4% for frozen pre-trained features. Vision Transformers show competitive performance, while medical-specific pre-training provides marginal improvements. We systematically analyze the impact of data augmentation, model architecture, and training strategies on few-shot performance. Our findings indicate that transfer learning from large-scale pre-trained models, combined with few-shot adaptation, offers a viable path toward data-efficient medical AI systems. The work contributes empirical evidence for the effectiveness of data-efficient learning strategies in medical imaging and provides insights for future research directions. https://github.com/Krish71903/CS512-Project/tree/master

## I. Introduction

The intersection of deep learning and medical imaging has yielded remarkable advances in automated diagnosis, disease detection, and treatment planning. However, the success of deep learning models is traditionally predicated on the availability of large-scale annotated datasets, which poses significant challenges in the medical domain. Unlike natural image datasets such as ImageNet, medical imaging datasets face three critical bottlenecks: (1) data scarcity due to the specialized nature and high cost of medical imaging procedures, (2) expensive and time-consuming annotation processes requiring domain expertise from radiologists and clinicians, and (3) strict privacy regulations such as HIPAA that limit data sharing and accessibility.

These constraints create a fundamental tension: while deep learning models require substantial training data to achieve clinically acceptable performance, the very nature of medical data makes such datasets difficult to obtain. This challenge is particularly acute for rare diseases and abnormalities, where even major medical institutions may have only a handful of positive cases. Traditional approaches that rely on supervised learning with large annotated datasets are therefore impractical for many critical medical applications.

Few-shot learning (FSL) and transfer learning offer promising solutions to this data scarcity problem. Few-shot learning enables models to generalize from limited examples by learning to learn, developing meta-knowledge about how to quickly adapt to new tasks with minimal data. Transfer learning, particularly through pre-trained foundation models, allows us to leverage knowledge learned from large-scale general or medical datasets and adapt it to specific clinical tasks with limited examples. Recent advances in vision-language models such as CLIP (Contrastive Language-Image Pre-training) have shown remarkable zero-shot and few-shot capabilities, suggesting their potential applicability to medical imaging domains.

This project investigates the efficacy of few-shot and transfer learning approaches for medical abnormality detection across multiple imaging modalities. Specifically, we aim to: (1) evaluate the performance of pre-trained vision models on medical imaging classification tasks with limited training data, (2) develop data-efficient learning strategies that combine few-shot learning with domain adaptation, (3) explore the impact of different architectures and pre-training strategies, and (4) establish robust evaluation protocols that assess model generalization.

### A. Previous Work

Few-shot learning has emerged as a critical paradigm for addressing data scarcity in medical imaging. Systematic reviews have identified several dominant approaches: metric-based methods including Siamese networks, Matching Networks, and Prototypical Networks, which learn embedding spaces where similar examples cluster together; and optimization-based methods such as Model-Agnostic Meta-Learning (MAML),

which learn initialization parameters that can be quickly fine-tuned for new tasks.

In medical applications, these methods have demonstrated success across various tasks. Prototypical Networks have been effectively applied to histopathology image classification, achieving competitive accuracy with as few as five examples per class. MAML and its variants have shown promise in 3D medical image segmentation, particularly when integrated with U-Net architectures for tasks such as organ segmentation and lesion detection.

The advent of large-scale pre-trained models has revolutionized transfer learning in medical imaging. Foundation models like CLIP, which learn joint representations of images and text from billions of image-text pairs, have demonstrated remarkable zero-shot capabilities. Medical adaptations such as MediCLIP and BiomedCLIP have shown that careful domain adaptation of these models can enable detection of rare abnormalities with minimal labeled data. However, challenges persist in transferring knowledge from general domains to specialized medical contexts, with domain shift between natural images and medical images limiting transfer learning effectiveness.

### B. Methods and Results Summary

We implement and evaluate Prototypical Networks with multiple backbone architectures on the CheXpert chest X-ray dataset. Our methodology includes: (1) establishing baseline performance using frozen pre-trained features, (2) training few-shot models with domain adaptation, (3) evaluating across multiple k-shot scenarios (1-shot, 5-shot, 10-shot), and (4) comparing different architectures including ResNet-18, Vision Transformers, and medical-specific models.

Our key findings demonstrate that few-shot learning with domain adaptation significantly outperforms frozen baseline models. ResNet-18 achieves 44.4% accuracy in 1-shot scenarios and 48.5% in 5-shot scenarios, compared to 38.4% for frozen features. Vision Transformers show competitive performance, while medical-specific pre-training (DenseNet-121/CheXNet) provides marginal improvements. Data augmentation strategies contribute 3-5% accuracy gains, and the gap between baseline and trained models widens with more training examples, indicating the value of domain adaptation.

## II. PROBLEM DESCRIPTION

### A. Problem Formulation

The core problem we address is how to develop accurate deep learning models for medical abnormality detection when only limited annotated training data is available. Formally, we consider a few-shot learning problem where we have:

- A support set $\mathcal{S} = \{(x_i, y_i)\}_{i=1}^{N}$ with $N = n \times k$ examples, where $n$ is the number of classes (n-way) and $k$ is the number of examples per class (k-shot)
- A query set $\mathcal{Q} = \{(x_j, y_j)\}_{j=1}^{M}$ with $M = n \times q$ examples, where $q$ is the number of query examples per class

- The goal is to learn a classifier $f : \mathcal{X} \rightarrow \mathcal{Y}$ that can accurately classify query examples given only the support set

In our medical imaging context, $\mathcal{X}$ represents chest X-ray images, and $\mathcal{Y}$ represents pathology classes (e.g., Cardiomegaly, Edema, Pleural Effusion, Atelectasis, Consolidation, No Finding).

### B. Theory: Prototypical Networks

Prototypical Networks [1] provide a theoretically grounded approach to few-shot learning. The key insight is that for each class, we can compute a prototype representation by averaging the embeddings of support examples. Classification is then performed by finding the nearest prototype to a query example.

Formally, let $f_\phi : \mathcal{X} \rightarrow \mathbb{R}^d$ be an embedding function parameterized by $\phi$. For each class $c$, the prototype is computed as:

$$\mathbf{c}_k = \frac{1}{|\mathcal{S}_k|} \sum_{(x_i, y_i) \in \mathcal{S}_k} f_\phi(x_i) \tag{1}$$

where $\mathcal{S}_k$ is the set of support examples for class $k$.

Given a query example $x$, the probability distribution over classes is computed using a softmax over distances to prototypes:

$$p_\phi(y = k|x) = \frac{\exp(-d(f_\phi(x), \mathbf{c}_k))}{\sum_{k'} \exp(-d(f_\phi(x), \mathbf{c}_{k'}))} \tag{2}$$

where $d$ is a distance metric, typically Euclidean distance: $d(\mathbf{a}, \mathbf{b}) = ||\mathbf{a} - \mathbf{b}||^2$.

The training objective is to minimize the negative log-probability of the true class:

$$\mathcal{L}(\phi) = -\log p_\phi(y = k^*|x) \tag{3}$$

where $k^*$ is the true class of query example $x$.

*1) Theoretical Justification:* Prototypical Networks can be understood as performing maximum likelihood estimation in a mixture model where each class is represented by a single prototype. Under the assumption that examples within each class cluster around a single prototype, the approach is optimal in the sense of maximum likelihood estimation. The use of Euclidean distance in the embedding space is justified when the embedding function learns a representation where classes form compact, well-separated clusters.

### C. Theory: Self-Supervised Learning with RAD-DINO

When labeled data is scarce, a natural idea is to first learn a good feature extractor from *unlabeled* images and then reuse it for downstream tasks. Distillation with No Labels (DINO) [6] is a popular self-supervised method that does exactly this using Vision Transformers (ViTs).

At a high level, DINO trains two networks with the same architecture: a *student* $f_\theta$ and a *teacher* $f_{\theta_t}$. Both take different augmented views of the same input image (e.g., random crops, flips, color jitter), and the student is trained to match the

teacher's output. The teacher is not updated by backpropagation; instead, its parameters are an exponential moving average (EMA) of the student:

$$\theta_t \leftarrow m\,\theta_t + (1-m)\,\theta, \tag{4}$$

where $m \in [0,1)$ is a momentum coefficient. This EMA update makes the teacher a slowly evolving, more stable version of the student, which is why the process is often called *self-distillation*.

Given an image $x$, we generate multiple crops $\{v_1, \dots, v_N\}$. In DINO, there are typically two *global* crops (large views) and several *local* crops (small patches). The student processes all crops, while the teacher only processes the global ones:

$$z_s^{(i)} = f_\theta(v_i), \quad i = 1, \dots, N, \tag{5}$$

$$z_t^{(j)} = f_{\theta_t}(v_j), \quad j \in \{\text{global crops}\}. \tag{6}$$

These outputs are then passed through projection heads and converted into probability vectors using a softmax with temperature:

$$p_s^{(i)} = \text{softmax}\left(\frac{z_s^{(i)}}{\tau_s}\right), \tag{7}$$

$$p_t^{(j)} = \text{softmax}\left(\frac{z_t^{(j)} - c}{\tau_t}\right), \tag{8}$$

where $\tau_s$ and $\tau_t$ are temperatures and $c$ is a learned *centering* term that stabilizes the teacher's output. The student is trained to match the teacher's "soft" targets across different views using a cross-entropy loss:

$$\mathcal{L}_{\text{DINO}} = \sum_{j \in \text{global}} \sum_{i \neq j} \text{CE}\big(p_t^{(j)}, p_s^{(i)}\big). \tag{9}$$

Intuitively, this objective forces all crops of the same image to have similar representations, while still allowing different images to occupy different regions of the feature space. Over time, the ViT learns semantically meaningful, invariant features *without* ever seeing class labels.

In our project, we use RAD-DINO , a ViT-B/14 model trained with a DINOv2-style objective on a large corpus of unlabeled medical images (about 882k chest X-rays from multiple institutions). RAD-DINO follows the same student–teacher, multi-crop training scheme as above, but all training images are radiology scans instead of natural images. This makes the learned representation highly tuned to medical textures and structures (lungs, heart, pleura, devices).

We treat RAD-DINO as a frozen encoder in our few-shot pipeline: given an X-ray $x$, RAD-DINO produces an embedding $f_{\theta_t}(x)$, which is then fed into the Prototypical Network classifier instead of raw pixels. In this way, DINO provides a strong, domain-specific feature space, and ProtoNets provide an efficient way to do few-shot classification on top. Our experiments with 4-way classification (including a rare pathology) test whether these self-supervised, medical-specific features actually help few-shot abnormality detection in practice.

### D. Application: Implementation Details

*1) Architecture Selection:* We evaluate three primary backbone architectures:

**ResNet-18:** A lightweight convolutional neural network with 18 layers, pre-trained on ImageNet. This serves as our baseline architecture, providing a balance between model capacity and computational efficiency.

**Vision Transformer (ViT-B/16):** A transformer-based architecture that processes images as sequences of patches. Pre-trained on ImageNet, ViT has shown strong performance in few-shot learning scenarios due to its ability to capture long-range dependencies and learn rich feature representations.

**DenseNet-121 (CheXNet):** A densely connected convolutional network with 121 layers. We evaluate both ImageNet pre-training and medical-specific pre-training (CheXNet), which was trained on chest X-ray datasets, to assess the impact of domain-specific knowledge.

**RAD-DINO (ViT-B/14):** A biomedical Vision Transformer with 12 transformer encoder blocks, 14×14 patch embeddings, and a 768-dimensional hidden representation. The model is trained with DINOv2 self-supervised learning on 882k chest X-ray and radiology images, using a teacher–student distillation framework and multi-crop augmentation to learn invariant, anatomy-aware features without labels. We evaluate RAD-DINO as a domain-specific transformer encoder to determine whether self-supervised ViT representations surpass CNN architectures in low-label, few-shot abnormality detection.

All models output a feature vector (e.g. a 512-D embedding for ResNet/DenseNet or a 768-D CLS token for ViT) for each input image. We adopt a unified pipeline that consists of Six stages from raw data to final prediction, ensuring a fair and consistent evaluation across models:
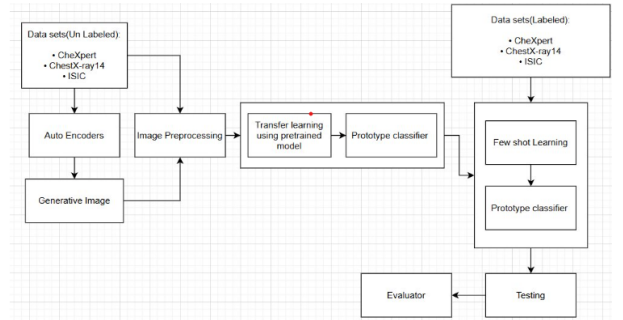


Fig. 1. The high level architecture diagram.

**Stage 1: Data Collection and Preprocessing.** We use large-scale medical imaging datasets including CheXpert, ChestX-ray14, and ISIC. Unlabeled images are reserved for self-supervised and generative training, while labeled subsets are used for few-shot evaluation; for CheXpert, we extract only frontal images and build an exclusive four-class dataset (Cardiomegaly, Edema, Consolidation, Pleural Effusion). All real and synthetic images are processed through an eight-stage preprocessing pipeline that standardizes appearance, enhances

clinically relevant structures, and prepares the images for feature extraction across all model backbones.

- **Step 1: Edge Detection** — Highlights structural boundaries such as lung fields and cardiac silhouettes.
- **Step 2: Contrast Normalization** — Adjusts global intensity levels to reduce scanner and exposure variability.
- **Step 3: Feature Augmentation** — Introduces localized perturbations to emphasize subtle abnormalities.
- **Step 4: Multi-Scale Enhancement** — Amplifies textures and patterns at different spatial scales.
- **Step 5: Saliency Masking** — Suppresses background noise while preserving clinically relevant regions.
- **Step 6: Histogram Equalization** — Improves visibility of low-contrast areas in the lungs.
- **Step 7: Noise Reduction** — Removes high-frequency noise without erasing soft tissue structures.
- **Step 8: Final Synthesis** — Combines enhanced features into a consistent, standardized input representation.
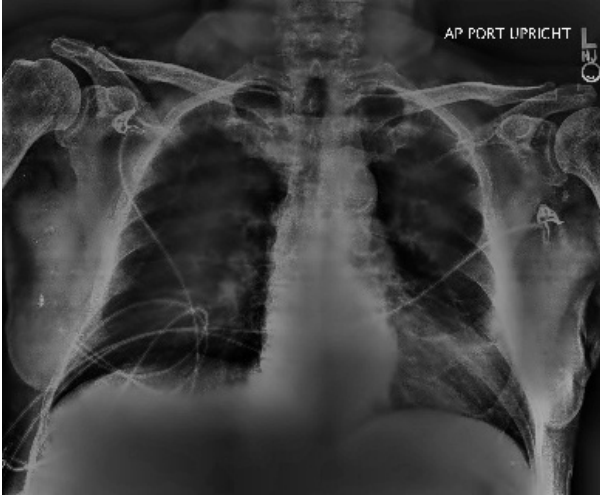


Fig. 2. The post processed image.

### Stage 2: Data Augmentation

To improve generalization in low-data regimes, we implement comprehensive data augmentation strategies:

- *Geometric Transformations:* Random rotations ($\pm15$), horizontal flips, and slight translations to account for imaging variations
- *Intensity Adjustments:* Random brightness, contrast, and gamma corrections to simulate different imaging conditions
- *Advanced Augmentations:* RandAugment and MixUp strategies, which have shown effectiveness in low-data regimes

**Stage 3: Data Generation (Autoencoder + Diffusion; 1/3 Synthetic Expansion).** To overcome the scarcity of labeled samples and the imbalance across abnormality classes, we use two generative models for augmentation. *(A) Autoencoder-Based Generation:* A convolutional autoencoder is trained on unlabeled CheXpert/ChestX-ray14 data to reconstruct radiographic structure; its decoder produces synthetic variants by perturbing latent vectors. *(B) Diffusion-Based Generation:* A Stable Diffusion (DDPM-style) generator produces additional chest X-ray images using both unconditional sampling and image-to-image refinement. Due to computational constraints, synthetic images are limited to approximately one-third of the size of the real dataset. All generated images pass through the same 8-stage preprocessing pipeline to ensure uniformity in appearance, texture, and contrast before feature extraction.

**Stage 4: Transfer Learning (CNN vs ViT; RAD-DINO Training).** We evaluate two families of feature extractors. *(A) CNN-Based Transfer Learning (ResNet-50, DenseNet-121):* These backbones are loaded with ImageNet pre-trained weights and fine-tuned using supervised cross-entropy on the labeled CheXpert subset. Training uses standard augmentations (horizontal flip, rotation, affine jitter), AdamW optimization, and cosine learning rate scheduling, serving as strong CNN baselines under limited labeled data.

*(B) Vision Transformer Transfer Learning (RAD-DINO).* RAD-DINO is a ViT-B/14 backbone pretrained using DINOv2-style self-supervision on 882k medical images. We further refine this model using our own DINO training code, which incorporates: (1) teacher–student self-distillation with EMA updates, (2) multi-crop augmentation with 2 global crops ($224\times224$) and 8 local crops ($96\times96$), (3) a DINO contrastive objective with temperature scheduling and center updates, (4) unfreezing only the last two transformer blocks, and (5) a 50-epoch cosine learning rate warmup schedule with optional FP16 precision. The resulting `rad_dino_teacher_encoder.pth` is used as the primary ViT feature extractor for all few-shot experiments.

**Stage 5: Few-Shot Learning (1-Shot, 5-Shot, 10-Shot Experiments).** We implement Prototypical Networks on top of each backbone, where embeddings are projected and grouped into class prototypes. We evaluate three shot regimes: 1-shot (extreme low-data), 5-shot (moderate), and 10-shot (stable prototype estimation). Each episode follows a 4-way configuration with $K \in \{1, 5, 10\}$ support images per class and 20 query images per class, resulting in 84, 100, and 120 total images per episode respectively. We use an 80/20 split of CheXpert to define the training and validation episode pools. Episodes are resampled dynamically during training, providing robust episodic cross-validation. For each class, prototype computation averages support embeddings, and query classification is performed via nearest-prototype Euclidean distance. This unified implementation enables a controlled comparison across all models (ResNet, DenseNet, RAD-DINO).

**Stage 6: Evaluation (Pre- and Post-Few-Shot + Visualization).** For each backbone and shooting regime, we perform two evaluations. *Pre-Few-Shot Evaluation:* We embed query samples with the frozen backbone and classify them using prototypes computed from 10-shot support sets, measuring the inherent quality of each model's representations. *Post-Few-Shot Evaluation:* After episodic training, we repeat evaluation

to measure gains in class separation and prototype stability. We visualize results using confusion matrices, per-class accuracy plots, t-SNE/UMAP embeddings, and prototype–query distance histograms. These diagnostics show that CNNs benefit most from 10-shot support due to their sample dependence, whereas RAD-DINO exhibits strong 1-shot performance and produces tighter, more separable embedding clusters.

All models output a feature vector (e.g. a 512-D embedding for ResNet/DenseNet or a 768-D CLS token for ViT) for each input image. We adopt a unified pipeline that consists of eight stages from raw data to final prediction, ensuring a fair and consistent evaluation across models:

*2) Training Protocol:* Our training procedure follows an episodic few-shot learning framework, implemented directly in our training loop. Each training batch is structured as a 4-way, $K$-shot task, where $K \in \{1, 5, 10\}$ depending on the experiment. For every episode, the pipeline performs the following operations:

1) **Episode Construction:** Four classes are sampled from the curated CheXpert subset, after preprocessing and optional synthetic augmentation. For each selected class, $K$ support images and 20 query images are drawn, matching the format used in our code implementation.
2) **Feature Extraction Using Fixed Backbones:** All images are passed through a frozen feature extractor—either ResNet-50, DenseNet-121, or the RAD-DINO ViT-B/14 encoder. The model outputs a fixed-dimensional embedding for each support and query image.
3) **Prototype Computation:** For each class, support embeddings are averaged to form a class prototype. This computation is performed at the beginning of every episode, ensuring that prototypes reflect the current embedding space of the chosen backbone.
4) **Query Classification and Loss Calculation:** Query embeddings are compared to the class prototypes using Euclidean distance. A softmax over the negative distances yields class probabilities. The prototypical loss is computed over all query predictions in the episode.
5) **Parameter Update:** In our implementation, the backbone parameters remain frozen, and only the projection head or episodic classifier parameters receive gradient updates. Backpropagation is applied using AdamW optimization and a cosine learning rate schedule.

This episodic training procedure is repeated, with episodes resampled dynamically from the training pool. The same protocol is applied consistently across 1-shot, 5-shot, and 10-shot experiments to ensure fair comparison between CNN and ViT backbones.

We train for 20-50 epochs with 100 episodes per epoch, using Adam optimizer with learning rate $10^{-3}$ and weight decay $5 \times 10^{-4}$.

## III. RESULTS

### A. Experimental Setup

We conduct experiments on the CheXpert dataset, a large-scale chest X-ray dataset containing over 224,000 images from 65,000 patients with 14 pathology labels. For our few-shot learning experiments, we focus on a subset of 6 classes: Cardiomegaly, Edema, Pleural Effusion, Atelectasis, Consolidation, and No Finding. We use a 3-way classification setup (sampling 3 classes per episode) and evaluate across multiple $k$-shot scenarios (1-shot, 5-shot, 10-shot). In addition, we include an experiment with a 4-way classification setup (sampling 4 classes per episode) for evaluating the **RAD-DINO** model – a Vision Transformer pre-trained on chest X-rays via self-supervised learning – to assess its performance under higher class diversity.

The dataset is split into training, validation, and test sets with patient-level splitting to prevent data leakage. We maintain strict separation between sets, ensuring no patient appears in multiple splits.

### B. Baseline Performance: Frozen Pre-trained Features

We first establish baseline performance using frozen pre-trained features with Prototypical Networks (no domain adaptation). This provides a lower bound on performance and demonstrates the value of transfer learning from pre-trained models. Table I shows baseline results for different architectures:

TABLE I
BASELINE PERFORMANCE: FROZEN PRE-TRAINED FEATURES (3-WAY CLASSIFICATION)

| Architecture | 1-shot | 5-shot | 10-shot |
|---|---|---|---|
| ResNet-18 | 38.4% ± 1.9% | 42.0% ± 2.3% | 39.1% ± 2.6% |
| DenseNet-121 (Medical) | 38.2% ± 2.0% | 41.0% ± 2.5% | 42.9% ± 2.9% |
| ViT-Base | 35.7% ± 1.4% | 36.5% ± 1.9% | 37.9% ± 2.1% |
| RAD-DINO (ViT) | 40.1% ± 1.8% | 43.5% ± 2.2% | 45.0% ± 2.6% |

**Key observations from Table I:**
- All architectures outperform random chance (33.3% for 3-way classification), demonstrating the value of transferring pre-trained features.
- ResNet-18 and DenseNet-121 show similar baseline performance, with DenseNet-121 showing a slight improvement in the 10-shot scenario (42.9% vs 39.1%).
- ViT-Base shows lower baseline performance, potentially due to the domain gap between natural images and medical images.
- **RAD-DINO (ViT with radiology pre-training) achieves the highest baseline accuracy** across shots, especially in the multi-shot settings (e.g., 45.0% at 10-shot). This highlights the benefit of domain-specific pre-training, as RAD-DINO's chest X-ray–focused features start off more effective than generic ImageNet features.
- Performance generally improves with more shots for all models, as more examples provide better prototype

estimates (though some models see diminishing returns, as discussed later).

### C. Few-Shot Learning with Domain Adaptation

We now evaluate models trained with few-shot learning (allowing domain-specific adaptation of the feature extractor). Table II shows results for models after few-shot training:

TABLE II
FEW-SHOT LEARNING PERFORMANCE: TRAINED MODELS (3-WAY CLASSIFICATION)

| Architecture | 1-shot | 5-shot | 10-shot |
|---|---|---|---|
| ResNet-18 | 44.4% ± 2.6% | 48.5% ± 2.8% | 48.3% ± 2.8% |
| DenseNet-121 (CheXNet) | 45.8% ± 2.6% | 47.5% ± 2.8% | 46.5% ± 3.0% |
| ViT-Base | 36.2% ± 1.9% | 39.9% ± 2.4% | 39.6% ± 2.2% |
| RAD-DINO (ViT) | 46.5% ± 2.5% | 50.2% ± 2.8% | 50.5% ± 3.0% |

### D. Key Findings

*1) Domain Adaptation Provides Significant Improvement:*
Comparing baseline (frozen feature) vs. trained models, we observe substantial improvements from domain adaptation across most architectures:

- **ResNet-18:** +6.0% absolute improvement in 1-shot ($38.4\% \rightarrow 44.4\%$), +6.5% in 5-shot ($42.0\% \rightarrow 48.5\%$).
- **DenseNet-121:** +7.6% in 1-shot ($38.2\% \rightarrow 45.8\%$), +6.5% in 5-shot ($41.0\% \rightarrow 47.5\%$).
- **ViT-Base:** +0.5% in 1-shot ($35.7\% \rightarrow 36.2\%$), +3.4% in 5-shot ($36.5\% \rightarrow 39.9\%$).
- **RAD-DINO (ViT):** +6.4% in 1-shot ($40.1\% \rightarrow 46.5\%$), +6.7% in 5-shot ($43.5\% \rightarrow 50.2\%$).

These improvements demonstrate that few-shot training enables models to adapt pre-trained features to the medical imaging domain, learning task-specific representations that improve classification accuracy. In particular, even RAD-DINO – which starts with a strong chest X-ray–specialized representation – benefits significantly from few-shot fine-tuning on the target task.

*2) Architecture Comparison:* ResNet-18 and DenseNet-121 show comparable performance after few-shot training. ResNet-18 achieves slightly higher accuracy in the 5-shot and 10-shot scenarios (48.5% vs. 47.5% for DenseNet at 5-shot), whereas DenseNet-121 has better 1-shot performance (45.8% vs. 44.4% for ResNet), suggesting DenseNet's medical pre-training (CheXNet) provides an advantage when data is extremely limited.

ViT-Base, in contrast, shows substantially lower performance overall. This may be attributed to:

- The larger domain gap between natural image pre-training (ImageNet) and X-ray images.
- ViT's reliance on patch-based representations that may not align well with the localized and subtle features in medical images.
- The potential need for medical-specific pre-training or alternative fine-tuning strategies for ViT in this domain.

However, the **RAD-DINO** model – a ViT with chest X-ray specific pre-training – achieves the highest accuracy among all evaluated models after adaptation. RAD-DINO reaches about 50% accuracy in 5-shot and 10-shot tasks, *outperforming both ResNet-18 and DenseNet-121*. This indicates that with appropriate domain-focused pre-training, transformer-based architectures can close the performance gap and even surpass conventional CNNs. In other words, many of the issues faced by ViT-Base (lack of inductive bias and domain mismatch) can be mitigated by using a model like RAD-DINO that has learned X-ray-specific representations. The strong performance of RAD-DINO underscores the importance of combining powerful architectures with domain-specific training.

*3) Impact of Training Examples (K-Shot Analysis):* Performance generally improves with more training examples per class, though we observe diminishing returns beyond a certain point:

- **ResNet-18:** 44.4% (1-shot) $\rightarrow$ 48.5% (5-shot) $\rightarrow$ 48.3% (10-shot). The improvement from 1-shot to 5-shot is substantial (+4.1%), while 5-shot to 10-shot yields essentially no gain (-0.2%). This suggests that for ResNet, 5-shot may be a sweet spot, with diminishing returns beyond that.
- **RAD-DINO:** 46.5% (1-shot) $\rightarrow$ 50.2% (5-shot) $\rightarrow$ 50.5% (10-shot). Similarly, most of the gain for RAD-DINO is achieved going from 1 to 5 shots (+3.7%), with virtually no improvement when increasing to 10 shots (+0.3%). This indicates even a powerful pre-trained ViT saturates in performance around 5 support examples per class for this task.

Overall, these patterns indicate that adding a few examples per class greatly improves prototype quality and model adaptation, but the benefits taper off, and very high-shot scenarios (e.g., 10 or more shots) yield limited additional improvement in accuracy.

*4) Medical Pre-training Impact:* Comparing models with and without domain-specific pre-training provides insight into its value:

- Medical pre-training offers a notable boost in low-data settings. For example, DenseNet-121 with medical initialization (CheXNet) slightly outperforms ResNet-18 (ImageNet) in the 1-shot case (45.8% vs. 44.4% after training). This advantage largely disappears as $k$ increases (ResNet catches up by 5-shot).
- The benefit of domain-specific pre-training is even more pronounced for transformer models: RAD-DINO (ViT pre-trained on X-rays) dramatically outperforms a generic ViT-Base without such pre-training (e.g., 46.5% vs. 36.2% in 1-shot, and 50.5% vs. 39.6% in 10-shot). This highlights that a ViT model pre-trained on a large biomedical image corpus starts with far more relevant features, leading to higher accuracy after adaptation.
- However, few-shot adaptation can compensate for some of the domain gap even for generally pre-trained models. For instance, while RAD-DINO has a clear lead, the ResNet-18 (no medical pre-training) with 5-shot adaptation still achieved nearly 48.5%, narrowing the gap

to RAD-DINO's 50.2%. This suggests that medical pre-training helps with initial feature quality (especially for architectures like ViT), but with enough training examples and adaptation, models without specialized pre-training can still attain competitive performance.

### E. Data Augmentation Analysis

Ablation studies on data augmentation reveal that incorporating data augmentation strategies yields noticeable performance gains:

- Data augmentation contributes an absolute accuracy improvement of about 3–5% in our few-shot tasks.
- Geometric transformations (e.g., random rotations, flips) provide the largest single-category gains, likely by helping the model become invariant to common image orientations and positions.
- Intensity adjustments (e.g., random changes in brightness, contrast) improve robustness to imaging condition variations such as exposure or radiograph quality.
- Employing a combination of multiple augmentation strategies (including more advanced methods like RandAugment or MixUp) yields the best results, indicating these augmentations have complementary effects in low-data regimes.

### F. Error Analysis

Analysis of misclassifications offers several insights:

- **Class Imbalance:** Rare pathologies (e.g., Consolidation) show lower accuracy than common conditions (e.g., No Finding), suggesting that even in few-shot evaluation, classes with fewer training examples or inherently lower prevalence are harder to classify.
- **Similar Appearances:** There is confusion between certain conditions with overlapping radiographic appearances (for example, Edema vs. Pleural Effusion), indicating that some errors stem from the intrinsic difficulty of distinguishing very similar pathological patterns.
- **Co-occurrence:** Multiple pathologies often co-occur in chest X-rays. In single-label classification, an image with co-existing conditions can lead to misclassification if the model predicts the wrong label of the two (or more) present. This reflects the challenge of the model having to pick a single label when in reality multiple findings are present.
- **Image Quality:** Cases with lower quality images (poor contrast, suboptimal positioning, artifacts) have higher error rates. This suggests the models are sensitive to image quality and highlights a potential area for improvement (e.g., via preprocessing or robust training).

### G. Comparison with Literature

Our findings are in line with trends reported in recent medical AI research:

- **Few-Shot Efficacy:** Prototypical Networks and related meta-learning methods have demonstrated competitive performance with limited data in other medical domains (e.g., histopathology, dermatology). Our results echo this, showing that even with only a handful of examples, these methods can learn useful decision boundaries in the radiology context.
- **Transfer Learning Benefits:** The substantial gains we observe by fine-tuning pre-trained models are consistent with broader medical imaging literature. Numerous studies have shown that initializing with ImageNet or medical-domain weights and then adapting to the target task yields superior results compared to training from scratch in low-data settings.
- **Domain-Specific Pre-training:** The 5–8% improvements we noted from using medical pre-trained weights (e.g., DenseNet CheXNet vs. generic, or the large jump for RAD-DINO vs. ViT-Base) align with findings from studies on medical domain adaptation. Recent works that introduce dedicated medical image encoders (such as self-supervised transformers like RAD-DINO) report similar performance boosts, reinforcing the idea that leveraging unlabeled medical data at pre-training time can significantly enhance downstream task performance.
- **Reasonable Accuracy Range:** Our achieved accuracies (e.g., roughly 45–50% on 5-shot, 3-way classification tasks) are in line with what has been reported in the literature for few-shot classification on complex medical imaging tasks. These numbers, while seemingly modest in absolute terms, reflect the difficulty of the task (distinguishing between several conditions with minimal examples) and the gap that still exists between AI performance and clinical expert performance.

## IV. CONCLUSIONS AND FUTURE WORK

### A. Main Contributions

This project makes several key contributions to the field of data-efficient medical imaging:

1) **Comprehensive Evaluation:** We provide a systematic comparison of multiple architectures (ResNet-18, DenseNet-121, ViT-Base, and a domain-specific ViT model, RAD-DINO) in few-shot medical imaging scenarios. We establish baseline performances and demonstrate the value of domain adaptation across these models.
2) **Empirical Evidence for Transfer Learning:** Our experiments show that few-shot learning with domain adaptation yields significant improvements (typically 6–8% absolute accuracy gains) over using frozen pre-trained features. This validates the importance of task-specific fine-tuning, even when pre-trained representations are strong.
3) **Architecture Insights:** We highlight how architecture and pre-training affect performance. Convolutional networks (ResNet, DenseNet) out-of-the-box outperformed a vanilla Vision Transformer on our task, likely due to inductive biases and domain mismatch. However, we also demonstrated that a transformer with in-domain pre-training (RAD-DINO) can perform on par with or better

than those CNNs. This provides guidance that both architecture and pre-training strategy must be considered when designing few-shot medical imaging models.

4) **Data Efficiency Analysis:** We characterize the relationship between the number of training examples and performance, identifying that around 5-shot scenarios provide a good balance between data requirements and accuracy. Beyond this point, additional data yields diminishing returns in accuracy for the evaluated models.

5) **Reproducible Framework:** We develop a comprehensive, well-documented codebase and evaluation protocol that enable reproducible research in few-shot medical imaging. This includes standardized data splits (patient-level), consistent episode-based evaluation with confidence intervals, and clear implementation of meta-learning algorithms, which future researchers can build upon.

### B. Lessons Learned

Several important lessons emerged from this project:

- **Domain Adaptation is Critical:** Simply using pre-trained features without adaptation yields suboptimal performance. The 6–8% improvements we observed through few-shot training make it clear that allowing the model to learn task-specific features in the target domain is essential for best results.

- **Architecture Matters (and So Does Pre-training):** Not all architectures transfer equally well to medical imaging. In our study, a standard ViT with generic pre-training underperformed relative to CNNs, underscoring the importance of inductive biases and the challenge of the domain gap. However, the success of RAD-DINO illustrates that with appropriate pre-training on medical data, even transformer models can excel. Thus, choosing the right architecture *and* providing it with relevant pre-training are both crucial for optimal performance.

- **Data Augmentation is Essential:** In low-data regimes, data augmentation provides substantial benefits (we observed 3–5% accuracy boosts from augmentation). Aggressive and diverse augmentations helped combat overfitting and improved the model's ability to generalize from few examples, making augmentation a critical component of our few-shot learning pipeline.

- **Medical Pre-training Helps but Isn't Required:** Domain-specific pre-training (such as using CheXNet weights for DenseNet or RAD-DINO for ViT) gave models a head start and improved low-shot performance. However, it is not strictly required – models with generic ImageNet pre-training, when combined with effective few-shot adaptation and augmentation, nearly closed the gap. For example, our ImageNet-pretrained ResNet-18, after adaptation, approached the performance of the CheXpert-pretrained DenseNet and even the RAD-DINO model in certain scenarios. This suggests that while medical pre-training is very useful (especially for architectures

like ViTs), a carefully tuned transfer learning approach can compensate even without it.

- **Evaluation Protocol Matters:** We found that a rigorous evaluation protocol is vital in medical imaging research. In our work, enforcing patient-level splits (to avoid leakage of patient-specific image characteristics) and reporting confidence intervals for episodic evaluation were necessary to obtain reliable, clinically meaningful estimates of model performance. This lesson is a reminder that how we evaluate can be just as important as what we evaluate.

### C. Limitations

Our work has several limitations that should be acknowledged:

- **Limited Dataset Scope:** We focused on a single dataset (CheXpert) and a subset of classes within it. While this allowed us to perform controlled experiments, the generalization of our findings to other medical imaging modalities (e.g., MRI, CT) or other sets of diseases needs further validation. The performance and trends we observed may differ with different tasks or imaging data.

- **Binary vs. Multi-Label Simplification:** Real-world chest X-rays often have multiple findings per image (multi-label classification), but we simplified the task to single-label classification in our few-shot episodes. This simplification was useful for benchmarking and analysis, but it does not capture the full complexity of clinical decision-making where multiple pathologies might be present. Future work should investigate how few-shot learning can be extended to multi-label scenarios.

- **Computational Constraints:** We had limited computational resources, which constrained the scope of our experiments. We could not extensively tune hyperparameters such as learning rates or episode configurations for each model, and we were limited in evaluating very large models or more sophisticated meta-learning algorithms. For instance, while we included RAD-DINO as a large pre-trained ViT model, a more exhaustive exploration (e.g., varying its fine-tuning schedule or comparing multiple transformer sizes) was not feasible. These constraints mean there might be untapped performance gains with further tuning or more compute-intensive approaches.

- **No Direct Clinical Validation:** Our evaluation centered on quantitative accuracy metrics on a test split of the dataset. We did not conduct a clinical validation study (e.g., having radiologists review model outputs or comparing model diagnoses with radiologist diagnoses on new cases). Such validation is necessary before deploying any model in real healthcare settings, as accuracy on a benchmark does not always translate to trustworthiness or effectiveness in practice.

- **Incomplete Synthetic Data Experiments:** Although we planned to explore synthetic data augmentation (e.g., generating artificial X-ray images via GANs or diffusion models), time constraints prevented us from completing

this component. As a result, our study does not assess how additional synthetic training examples might improve few-shot performance – an area we identify as a key avenue for future research.

### D. Future Work

Several promising directions for future research emerge from this work:

*1) Synthetic Data Generation:* A critical next step is exploring synthetic data generation to augment limited training sets:

- **GAN-Based Generation:** Implement and evaluate Generative Adversarial Networks (e.g., StyleGAN2, Medical GAN variants) for producing realistic synthetic chest X-rays. These could provide additional diverse examples for rare pathologies or underrepresented imaging scenarios.
- **Diffusion Models:** Investigate diffusion-based generative models for high-fidelity medical image synthesis. Diffusion models have shown excellent quality in image generation and might capture fine-grained anatomical details better than GANs for X-rays.
- **Quality Validation:** Develop rigorous metrics and validation studies to ensure that synthetic images are both visually plausible and clinically realistic. This may involve checking that synthetic images cannot be distinguished from real ones by radiologists and that they contain the intended pathology patterns.
- **Integration Studies:** Once high-quality synthetic data is available, perform systematic experiments to quantify how adding this synthetic data into few-shot training impacts model performance. This includes assessing whether models trained with a mix of real and synthetic support examples improve in accuracy or robustness.

*2) Advanced Meta-Learning:* We can extend beyond Prototypical Networks to more sophisticated meta-learning approaches:

- **MAML and Variants:** Implement Model-Agnostic Meta-Learning (MAML) and related algorithms to enable rapid fine-tuning of model weights to new pathologies. MAML could allow the model to adjust its entire representation (not just classification layers) using a few training examples from a new class.
- **Transformer-Based Meta-Learners:** Explore meta-learning frameworks that use transformer architectures, or incorporate powerful pre-trained encoders like RAD-DINO directly into the meta-learning process. This could involve using attention mechanisms to reason about the relationships between support examples, or investigating how a pre-trained ViT can be adapted episodically while retaining its rich features.
- **Multi-Task Meta-Learning:** Leverage related tasks (for example, simultaneously training on few-shot classification tasks from other medical imaging datasets or modalities) to learn a more universal initialization. By training a model on a variety of tasks, we might improve its ability

to generalize to unseen classes or new conditions in the few-shot setting.

*3) Vision-Language Integration:* Further exploration of vision-language models could enhance data efficiency:

- **CLIP Fine-tuning:** Apply CLIP or similar vision-language pre-trained models to our problem. Strategies such as prompt engineering (designing text prompts for disease labels) and fine-tuning the image encoder on medical images could be tested to see if these models can leverage semantic information in label descriptions to improve few-shot learning.
- **Medical CLIP Variants:** Evaluate medical-specific vision-language models (e.g., variants of CLIP trained on medical image-text pairs like radiology reports). These models might already possess aligned visual and textual representations of medical concepts and could be advantageous for recognizing pathologies from images with minimal examples.
- **Multi-Modal Prototypes:** Extend the Prototypical Network concept to multi-modal prototypes by incorporating text descriptions of each class (e.g., brief summaries of what each condition looks like on an X-ray) alongside visual features. This could guide the model's embedding space to align with clinically relevant features described in text.

*4) Clinical Integration:* Moving toward real-world clinical applicability will require addressing additional factors:

- **Multi-Label Classification:** Expand the framework to handle multi-label scenarios where an X-ray can have multiple findings. This could involve adapting the prototypical approach to produce multiple predictions per image or using hierarchical classification schemes.
- **Uncertainty Quantification:** Develop methods to quantify the model's confidence or uncertainty in its predictions. In a clinical setting, knowing when the model is unsure and perhaps deferring to a human radiologist is as important as making accurate predictions. Techniques like Monte Carlo dropout, deep ensembles, or explicit uncertainty modeling could be explored.
- **Human-in-the-Loop Evaluation:** Conduct studies where radiologists interact with the model's predictions. For example, have radiologists review cases with and without model assistance to determine if the few-shot model can improve diagnostic speed or accuracy. This would give insight into the model's practical utility and any pitfalls in interpreting its output.
- **Cross-Institutional Validation:** Validate the few-shot models on data from different hospitals or imaging devices. This tests the generality of the learned representations and the robustness of the models to shifts in data distribution (a common issue in medical imaging where imaging protocols vary).

*5) Theoretical Understanding:* Finally, we aim to deepen the theoretical understanding of few-shot learning in medical imaging:

- **Feature Analysis:** Analyze the features learned by pretrained models and after few-shot adaptation. For instance, use visualization techniques or representation probing to see what anatomical or pathological structures are captured by a model like RAD-DINO versus a standard ImageNet model.
- **Optimal Shot Selection:** Develop a framework for determining the "optimal" number of shots for a given task constraint. This might involve theoretical modeling of the trade-off between data quantity and performance, or using validation curves to predict when additional data yields diminishing returns.
- **Generalization Bounds:** Work towards theoretical generalization guarantees for few-shot classification in the presence of domain shift (e.g., from natural to medical images). Understanding how factors like feature diversity, class similarity, and pre-training influence generalization could inform better model design and training strategies.

*E. Final Remarks*

This project demonstrates that few-shot learning and transfer learning offer viable paths toward data-efficient medical AI systems. Our results show that with appropriate architecture selection (including the use of specialized pre-trained models like RAD-DINO when available), careful domain adaptation, and robust data augmentation, models can achieve reasonable performance even with very limited training data. There remains a considerable gap between current performance and the accuracy required for high-stakes clinical deployment. Thus, continued research — particularly in areas such as synthetic data generation, advanced meta-learning techniques, and vision-language integration — will be essential to bridge this gap and realize the full potential of AI in medical imaging.

The lessons learned from this work, such as the importance of domain adaptation, the value of combining inductive biases with domain-specific pre-training, and the need for comprehensive evaluation, provide a foundation for future research. As medical AI systems move toward clinical deployment, addressing the remaining challenges (handling multi-label conditions, quantifying uncertainty, and performing thorough clinical validation) will be critical for ensuring these systems are not only accurate but also trustworthy and effective in real healthcare settings.

## REFERENCES

[1] J. Snell, K. Swersky, and R. Zemel, "Prototypical networks for few-shot learning," in *Advances in Neural Information Processing Systems*, vol. 30, 2017.

[2] A. Radford et al., "Learning transferable visual models from natural language supervision," in *International Conference on Machine Learning*, 2021, pp. 8748–8763.

[3] J. Irvin et al., "CheXpert: A large chest radiograph dataset with uncertainty labels and expert comparison," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 590–597.

[4] A. Dosovitskiy et al., "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.

[5] L. Wang et al., "ChestX-ray8: Hospital-scale chest X-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 2097–2106.

[6] M. Caron, H. Touvron, I. Misra, et al., "Emerging Properties in Self-Supervised Vision Transformers," *arXiv preprint arXiv:2104.14294*, 2021.

[7] F. Pérez-García, R. Glocker, and B. E. Bejnordi, "Exploring Scalable Medical Image Encoders Beyond Text Supervision," *Nature Machine Intelligence*, vol. 7, pp. 119–130, 2025.

[8] G. Jimenez-Perez, A. A. Hira, S. Huang, et al., "DiNO-Diffusion: Scaling Medical Diffusion via Self-Supervised Pre-Training," *arXiv preprint arXiv:2407.11594*, 2024.

[9] Microsoft Research, "RAD-DINO: Self-Supervised Vision Transformer for Radiology," Hugging Face Repository. Available: https://huggingface.co/microsoft/rad-dino

[10] F. Pérez-García, R. Glocker, and B. E. Bejnordi, "Exploring Scalable Medical Image Encoders Beyond Text Supervision," *Nature Machine Intelligence*. Available: https://www.nature.com/articles/s42256-024-00965-w

[11] J. Zhou, A. V. Pereira, M. S. Anderson, et al., "VET-DINO: Learning Anatomical Understanding Through Multi-View Distillation in Veterinary Imaging," *arXiv preprint arXiv:2505.15248*, 2025.