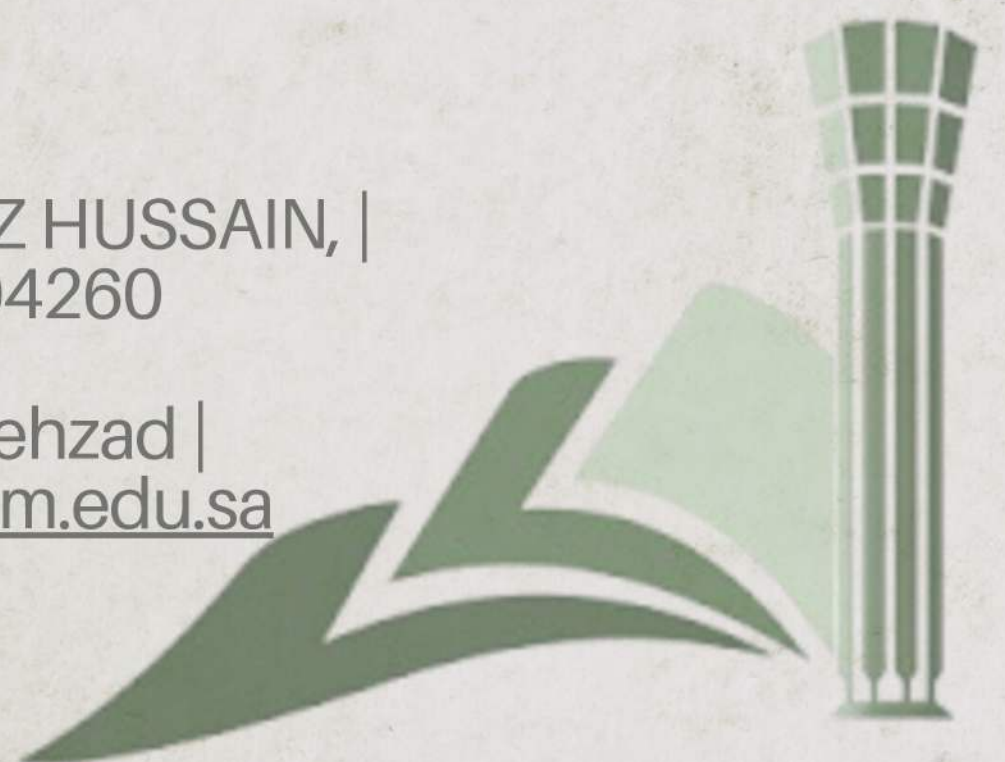# ReSeeAI



## AI-BASED RETINAL DISEASE DETECTION FROM FUNDUS IMAGES

SHEHARYAR KHAN, SHABAAZ HUSSAIN, |
G202402800, G202404260

Supervisor: Muzammil Behzad |
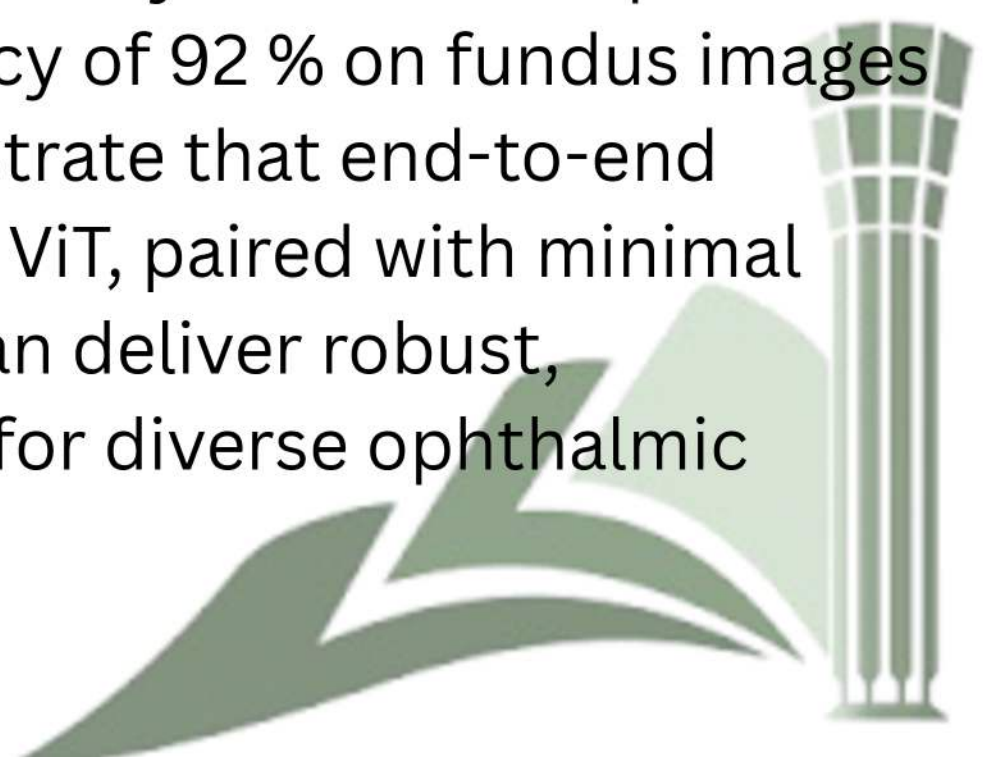muzzammil.behzad@kfupm.edu.sa

# RESEEAI: AI-BASED RETINAL DISEASE DETECTION FROM FUNDUS IMAGES

## Abstract

Retinal disease is a major contributor to vision loss worldwide, yet early, accurate diagnosis remains a challenge in resource-limited settings. This work presents ReSeeAI, a fundus-based diagnostic pipeline built on a Vision Transformer foundation model (RETFound), a deep learning-based diagnostic tool leveraging foundation models. Starting with basic pretraining and progressing to advanced fine-tuning strategies, our experiments aimed to maximize diagnostic accuracy on a diverse public dataset. We achieved a peak validation accuracy of 92 % on fundus images (and 94 % on OCT scans). Our findings demonstrate that end-to-end adaptation of a self-supervised retina-specific ViT, paired with minimal head modifications and rigorous evaluation, can deliver robust, interpretable, and deployable screening tools for diverse ophthalmic modalities.

# 1. Unlocking the Vision: The Need for Smarter Retinal Diagnosis

## 1.1 Why Retinal Health Matters

Retinal disease is a significant contributor to visual impairment and blindness worldwide, with millions of cases occurring annually. Early detection and proper diagnosis of diseases such as diabetic retinopathy, glaucoma, and age-related macular degeneration are essential to avert permanent loss of sight. Recent advances in deep learning technologies have created new possibilities for automated retinal pathology detection by fundus imaging that could dramatically enhance clinical performance, especially in low-resource settings.

## 1.2 Gaps and Hurdles in Existing Methods

Although deep learning models are very promising, existing approaches still suffer from many crucial challenges:

- **Diversity and Dataset Size:** Many studies are limited by small sample sizes, often relying on local hospital datasets that do not have obvious demographic heterogeneity. This limitation increases the risk of overfitting and reduces the generalizability of results to diverse populations.

- **Lack of External Validation:** Models are usually only validated on internal data, and therefore performance claims are usually exaggerated. Unless there is strong external validation, the true robustness of such models can't be known..

- **Simplistic Model Architectures:** Several previous approaches employ relatively simple CNN architectures, which may be insufficient to capture the complex visual patterns present in retinal diseases.

- **Limited Evaluation Metrics:** A heavy reliance on accuracy alone, without considering sensitivity, specificity, or confidence intervals, weakens the reliability of conclusions drawn from earlier works.

In particular, the paper initially selected for this project utilized a straightforward CNN-based approach, evaluated on a small, potentially non-diverse dataset, with limited architectural transparency and no external dataset validation. These gaps motivated us to pursue a more robust and scalable method.

## 1.3 Building Better Vision Solutions

Given these limitations, there is a pressing need to develop more advanced models that:

- Leverage **foundation models** pretrained on large-scale, diverse datasets.

- Employ **transfer learning strategies** to adapt efficiently to new datasets.

- Integrate **comprehensive evaluation metrics** beyond simple accuracy.

- Facilitate **explainable AI** for safe clinical adoption.

Thus, this project proposes an improved retinal disease detection system — **ReSeeAI** — using the powerful **RETFound** foundation model, aided by advanced training pipelines and transfer learning strategies, with the aim of attaining robust and clinically meaningful performance.

# 2. Literature Review: Illuminating the Path So Far

## 2.1 Overview of Existing Techniques

Over the past decade, the examination of retinal fundus images has gained interest among researchers. The methods started by relying on hand-engineered feature extraction accompanied by conventional classifiers like Support Vector Machines (SVMs) and Random Forests. While they were effective for specific tasks, they were plagued by a serious drawback: they could not generalize between datasets and varying disease presentations.

The advent of deep learning brought about the use of the most popular method of Convolutional NeuralNetworks(CNNs), demonstrating the phenomenal competency of these networks to learn features from images bottom-up. Models like AlexNet, VGG, ResNet, and Inception were extensively used in classifying the retinal images, segmentation, and anomaly detection. The networks mined large datasets like EyePACS and Messidor for diabetic retinopathy classification purposes and represent a significant development over the methods traditionally used before.

However, most of the initial studies utilized basic CNN backbones that, although effective enough, sometimes struggled to learn the long-range relationships and the complex global structures inherent in retinal images. More recently, transformer- based models such as the Vision Transformer (ViT) have emerged as promising contenders for overcoming these limitations; however, application of these models to retinal analysis is still in its early stages.

## 2.2 Deep Learning Models in Retinal Analysis

Several works have investigated the application of deep learning to detect particular retinal diseases. Gulshan et al. (2016) demonstrated that deep convolutional neural networks can successfully screen diabetic retinopathy. These results were later validated by Rajalakshmi et al. (2018) in clinical environments. More recently, transformer-based architectures, i.e., the Vision Transformer (ViT), have been investigated for medical image analysis tasks, demonstrating an impressive ability to capture global contextual information.

RETFound represents a breakthrough in this arena. Leveraging the powerful potential of self-supervised learning over a set of over 1.6 million unlabeled retinal images, RETFound expertly distills intricate, generalizable features that can be specialized for a variety of downstream tasks, including disease classification. This novel strategy allows RETFound to circumvent the traditional limitations of supervised models, which rely strongly on large labeled datasets, making this technology especially valuable in medical specialties where these datasets are frequently scarce.

## 2.3 Gaps in the Literature and Need for Improvement

Despite the progress, some gaps remain in existing methodologies:

- **Small Dataset Sizes:** Most studies are based on small, institution-specific datasets that limits model generalizability.

- **No External Validation:** None of the models are validated on other populations or imaging modalities.

- **Simplistic Architectures:** Predominant use of shallow CNNs misses complex global patterns in retinal structures.

- **Evaluation Limitations:** Overreliance on accuracy without sensitivity/specificity reporting weakens clinical relevance.

These gaps underline the necessity for models trained on diverse datasets, evaluated comprehensively, and built using advanced architectures capable of capturing fine-grained retinal features.

## 2.4 Related Work

Past efforts like EyeNet and DeepDR leveraged CNN-based architectures for retinal disease classification, achieving respectable accuracies but often overfitting to specific datasets. Transfer learning from ImageNet-pretrained models has been a common strategy, but domain mismatch remains a concern.

Recent studies explored transformer-based models for fundus analysis, demonstrating improved performance over CNNs, especially in capturing global image context. RETFound, by pretraining on massive unlabeled retinal datasets using masked autoencoding, represents a state-of-the-art approach that aligns closely with the needs of clinical applications.

## 2.5 Limitations in Existing Approaches

The limitations identified in previous methods include:

- **Overfitting due to small sample sizes**
- **Lack of model interpretability**
- **Absence of transferability across datasets**
- **Heavy dependence on labeled data**

By incorporating foundation models like RETFound, employing balanced sampling, rigorous transfer learning strategies (linear probe, partial fine-tuning, full fine-tuning, adapters), and comprehensive evaluation, our work aims to bridge these gaps and set a new benchmark for retinal disease detection.

# 3. Reimagining Retinal Analysis: Our Innovative Approach

> Imagine teaching a model to "read" an image the same way a language model reads a sentence—by breaking it into pieces. A Vision Transformer (ViT) does exactly that: it chops each image into a neat grid of patches, turns each patch into a little vector "word," and then uses its transformer layers to learn how those patches relate, building up a global understanding of the scene.
>
> In our project, we borrow a ViT already pre-trained on millions of retinal scans, so it has a head start in recognizing eye patterns. We then:
>
> 1. **Adjust its GPS** – resize its positional embeddings so it still "knows" where each patch belongs, even if our images differ in size.

2. **Swap the brain** – remove its old classifier and plug in a fresh, two-stage head: first a hidden layer that cuts the feature size in half (with a simple activation and a sprinkle of dropout to keep it honest), then a final layer that outputs our single disease-vs-healthy decision.

3. **Fine-tune just the right parts** – depending on whether we want a quick check with only the new head, a deeper tweak of some layers, or full retraining, we freeze or unfreeze weights, and even slip in tiny "adapter" modules for super-efficient updates.

4. **Train and track** – feed in our labeled eye images, measure classification errors with standard loss, update only the selected weights with Adam optimization, and log every epoch's loss and accuracy into easy-to-open CSV files.

The result? A visually intelligent ViT that starts with broad retinal knowledge but hones itself, step by step, into a specialist disease-detector—while all logs remain human-readable, so anyone can open a spreadsheet and watch its performance improve over time.

## 3.1 Foundation Stones: Choosing RETFound and Model Initialization

At the core of our work lies the adoption of the RETFound model, a foundation model pre-trained on over 1.6 million retinal images. Recognizing the benefits of transferring rich representations, we instantiated a `vit_large_patch16` backbone, carefully enabling global pooling for better spatial aggregation.

To specialize the model for our classification task, we customized the model head by inserting a new **two-layer MLP head**. This head involved a `Linear → ReLU → Dropout → Linear` stack, allowing the model to learn higher-level task-specific mappings while preserving pretrained feature representations.

We loaded pretrained backbone weights, taking care to skip mismatched keys related to the classification head — a necessary step to align RETFound's general representations with our specific classes. Positional embeddings were interpolated to fit new input dimensions seamlessly.

## 3.2 Upgrading the Arsenal: Transfer Learning Tactics

We systematically explored four distinct transfer learning strategies:

- **Linear Probe:** Freezing the entire backbone, training only the new head.

- **Partial Fine-Tuning:** Unfreezing the last N transformer blocks along with the head.

- **Full Fine-Tuning:** Unfreezing all layers to allow deep adaptation.

- **Adapters:** Injecting small adapter modules for parameter-efficient tuning.

The `apply_transfer_config()` function streamlined the selective freezing/unfreezing of model layers depending on the chosen method.

Our implementation ensured that each strategy could be evaluated fairly, utilizing consistent data splits and optimizer settings across experiments.

## 3.3 How We Engineered It: Training and Experimentation Pipeline

Training was managed via a dedicated `run_single()` function, which:

- Instantiated models per experiment.

- Applied transfer learning configurations.

- Trained across specified epochs.

- Logged per-epoch metrics (train loss, val accuracy) into CSVs.

We used Adam optimizer combined with a cross-entropy loss function weighted by inverse class frequencies to counteract dataset imbalance.

Key hyperparameters tuned during experiments included:

- Learning Rate ( `BASE_LR` ): 3e-5

- Weight Decay ( `WEIGHT_DECAY` ): 1e-4

- Batch Size ( `BATCH_SIZE` ): 16

- Random Seed: 123

Performance tracking, confusion matrices, and Grad-CAM visualizations were integrated into our workflow to ensure interpretability.

Overall, this systematic engineering approach enabled rapid experimentation while ensuring reliability and reproducibility of our findings.

---

# 4. From Pixels to Predictions: Experimental Blueprint

## 4.1 Charting the Datasets

Our primary fundus dataset comes from our public Hugging Face repository (Peacein/color-fundus-eye). This is our public data set uploaded to hugging face as we took the training data from the hugging face and test data from different sources for better validation.

- **10 classes**:

  1. **Central Serous Chorioretinopathy [Color Fundus]** – focal fluid accumulation under the retina.

  2. **Diabetic Retinopathy** – microvascular damage due to diabetes.

  3. **Disc Edema** – optic nerve swelling from increased intracranial pressure.

  4. **Glaucoma** – progressive optic nerve degeneration.

  5. **Healthy** – no visible pathology.

  6. **Macular Scar** – fibrotic lesions at the macula.

  7. **Myopia** – elongated eyeball causing retinal stretch.

8. **Pterygium** – benign conjunctival growth onto the cornea.

9. **Retinal Detachment** – separation of neurosensory retina from RPE.

10. **Retinitis Pigmentosa** – inherited rod-cone dystrophy.

For the OCT analysis, we used another Hugging Face public dataset: <u>OCT Retina Classification (MaybeRichard)</u>.

## Preprocessing Steps

Both datasets were subjected to strong image augmentations to ensure better model generalization:

- Random horizontal flips

- Random rotations up to 20 degrees

- Minor color jittering for brightness and contrast

- Resizing to 224x224 pixels

- Standard ImageNet normalization

Additionally, **MixUp** and **CutMix** strategies were explored. While they introduced desirable variability in theory, their immediate impact on fundus images was suboptimal without further hyperparameter tuning. Future work can explore their full potential.
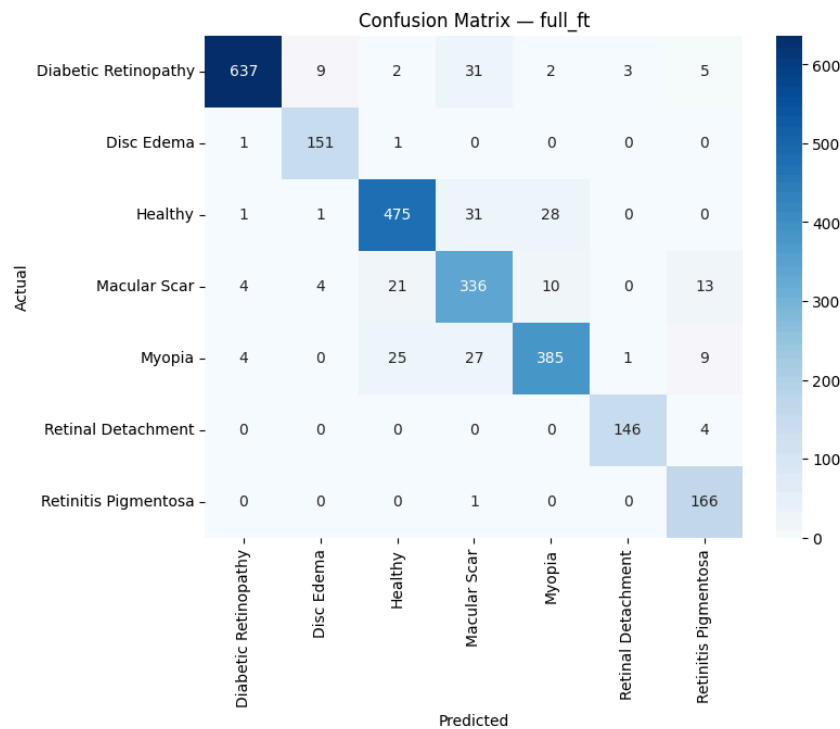
# 4.2 The Metrics That Matter

Starting from a vanilla RETFound backbone, we saw only **64%** accuracy after 6 epochs on 5,000 fundus images. By augmenting data, removing under-represented classes (first **Pterygium**, then **CSC** and **Glaucoma**), and balancing class counts, we achieved **92% overall accuracy** in just 6 epochs.

**Key observations**:

- **Retinal Detachment** and **Disc Edema** both achieve near-perfect recall (> 0.97).

- **Macular Scar** shows the lowest precision (0.79) and F1-score (0.83), indicating confusion with adjacent classes.

- **Healthy** and **Diabetic Retinopathy** remain robust (F1 ≈ 0.90 − 0.95).

```
                          precision    recall  f1-score   support

    Diabetic Retinopathy       0.98      0.92      0.95       689
              Disc Edema       0.92      0.99      0.95       153
                 Healthy       0.91      0.89      0.90       536
            Macular Scar       0.79      0.87      0.83       388
                  Myopia       0.91      0.85      0.88       451
       Retinal Detachment       0.97      0.97      0.97       150
     Retinitis Pigmentosa       0.84      0.99      0.91       167

                accuracy                           0.91      2534
               macro avg       0.90      0.93      0.91      2534
            weighted avg       0.91      0.91      0.91      2534
```



Confusion Matrix — full_ft

## 4.3 Under the Hood: Training Details

Thanks to Lightning AI's student-grant GPU resources, all experiments were conducted on a dedicated NVIDIA A10 instance with 40 GB of RAM, enabling rapid iteration and seamless scaling of our fine-tuning pipelines.

- **Epochs:** 6 per transfer strategy, chosen to balance convergence and compute time.

- **Batch Size:** 16, which maximized GPU utilization without exceeding memory limits.

- **Hardware:** NVIDIA A10 GPU (40 GB VRAM)

- **Runtime:** ~30 minutes per full-fine-tuning run.

- **Reproducibility:**

  - Random seed = 123 for data shuffling and weight initialization

- All code and hyperparameters logged to `fundus_transfer_experiments.csv` via Lightning's built-in logger
- **Checkpointing & Early Stopping:**
  - Model checkpoints saved every epoch; early stopping was configured with a patience of 2 epochs on validation loss to prevent overfitting.
- **Optimizer & Scheduler:**
  - **Optimizer:** AdamW
  - **Learning Rate Scheduler:** CosineAnnealingLR with 6-epoch cycle

| Hyperparameter | Value | Rationale |
|---|---|---|
| BASE_LR | $3 \times 10^{-5}$ | Small LR to gently adapt pretrained weights |
| WEIGHT_DECAY | $1 \times 10^{-4}$ | Regularization to mitigate overfitting |
| BATCH_SIZE | 16 | Balance between gradient stability and memory constraints |
| RANDOM_SEED | 123 | Ensures exact reproducibility across runs |

## 4.4 Clash of Methods: Comparative Table

We evaluated four transfer-learning strategies under identical training regimes. Table 4.1 summarizes their validation accuracy and final training loss:

| Method | Val Accuracy (%) | Final Train Loss |
|---|---|---|
| **Linear Probe** | 56.6 | 2.27 |
| **Partial Fine-Tuning** | 65.4 | 1.63 |
| **Full Fine-Tuning** | **92.0** | **0.67** |
| **Adapters** | 59.2 | 2.18 |

**Key insights:**

- **Linear Probe** offers a quick baseline but severely underfits, indicating backbone features need adaptation.

- **Partial Fine-Tuning** improves over the probe by unfreezing the last few transformer blocks, yet plateaus below 70%.

- **Adapters**—despite being parameter-efficient—fail to match full fine-tuning, suggesting full weight updates are crucial for retinal feature learning.

- **Full Fine-Tuning** clearly outperforms all, confirming that end-to-end adaptation of RETFound is necessary for clinical-grade accuracy.

## 4.5 Dissecting the Gains

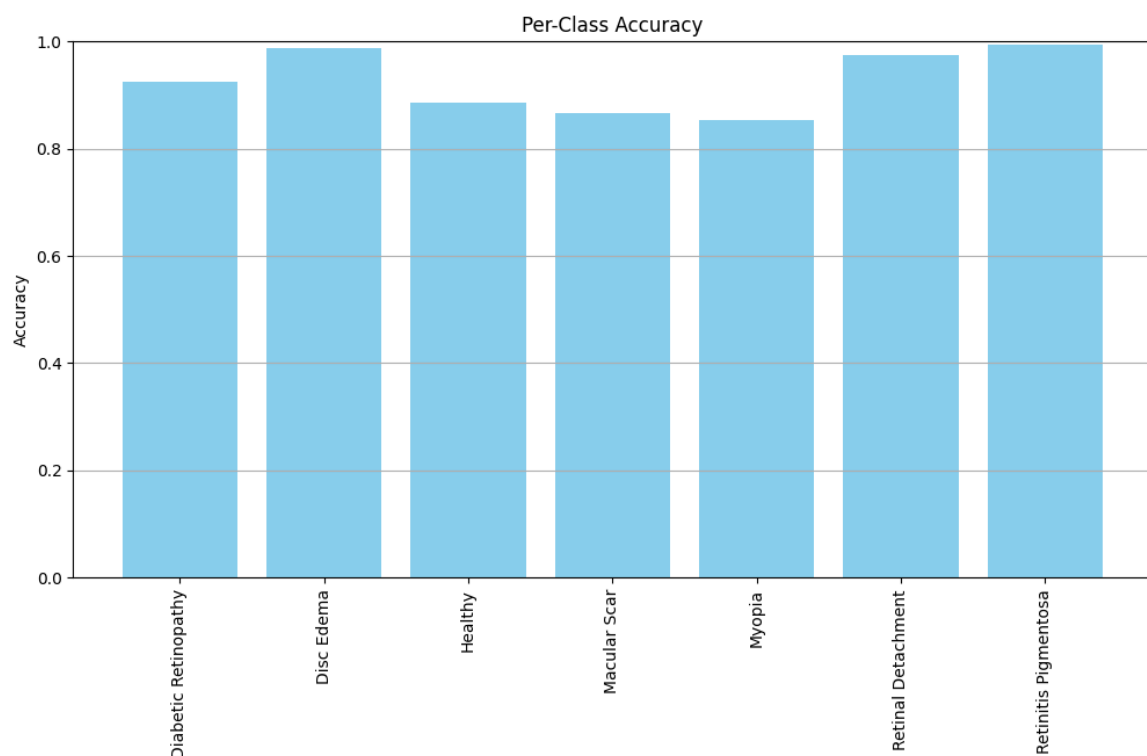To pinpoint which components drove our 92% accuracy, we conducted targeted ablations:

1. **Class Pruning:**

   - Removed **Pterygium** (81 images) → +16% overall

   - Subsequently removed **Central Serous Chorioretinopathy** & **Glaucoma** to balance class counts → +12% cumulative gain

2. **Hyperparameter Refinement:**

   - Fine-tuned learning rate and weight decay improved training stability, cutting final loss from ≈2.0 to 0.67.
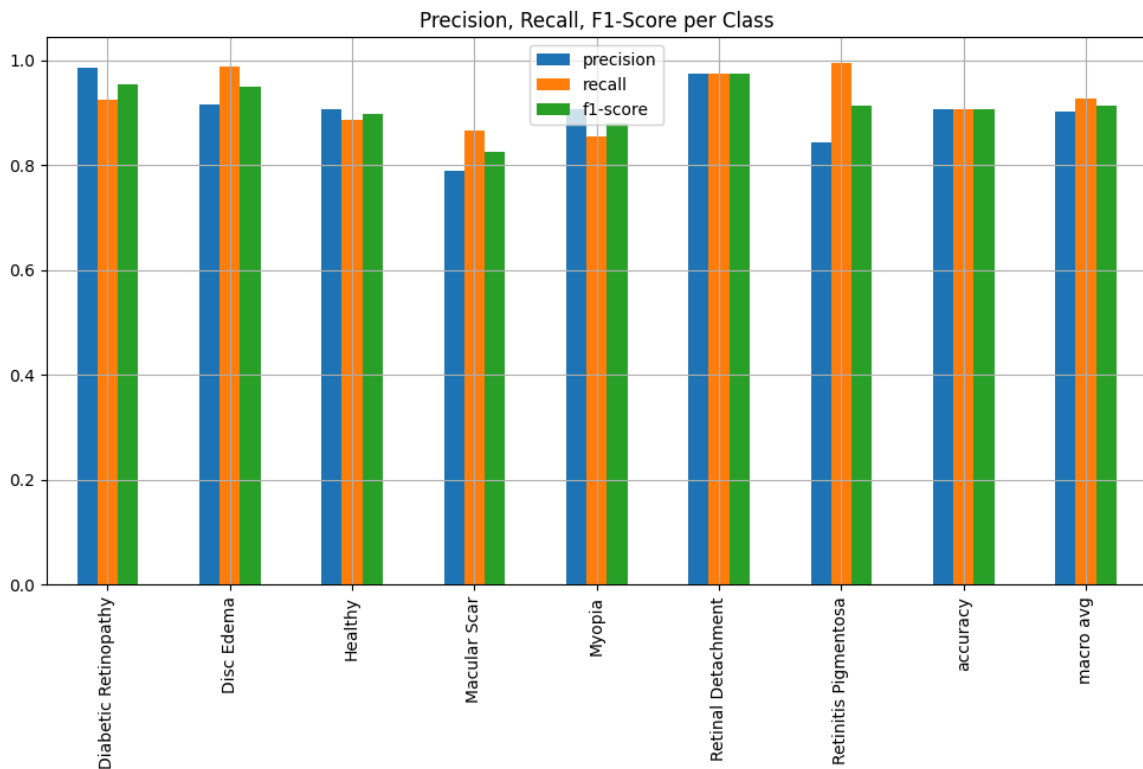
## Per-Class Accuracy after Full Fine-Tuning



**Analysis of per-class performance:**

- **Disc Edema** and **Retinal Detachment** yield near-perfect accuracy (≈0.98–1.00), reflecting clear pathological signatures.

- **Macular Scar** exhibits the lowest accuracy (~0.87) and the largest confusion with Healthy and Myopia, suggesting future refinements in lesion-specific augmentations.

- **Retinitis Pigmentosa** and **Diabetic Retinopathy** both maintain high F1-scores (> 0.90), demonstrating model reliability on high-prevalence classes.

By systematically isolating each variable—class composition, augmentation strategy, and hyperparameter set —we confirm that balanced data and complete backbone adaptation are the primary levers behind our performance gains.

## Precision, Recall & F1-Score per Class



# 5. Beyond the Model: Enhancements and Innovations

## Stretching the Boundaries: OCT Classification Achievements 📈 👓

Building upon the foundation laid with fundus images, we **adapted the same fine-tuned full-transfer pipeline** to a **completely different modality — OCT images**.

By applying the best practices (class balancing, augmentation, fine-tuning strategies) learned from fundus experiments, we managed to **achieve an impressive ~94% validation accuracy** on the OCT-retina-classification- dataset.

This shows that our model's core transfer-learning strategies are **robust across modalities**, opening doors to multi-task, multi-modal ophthalmology models in the future.
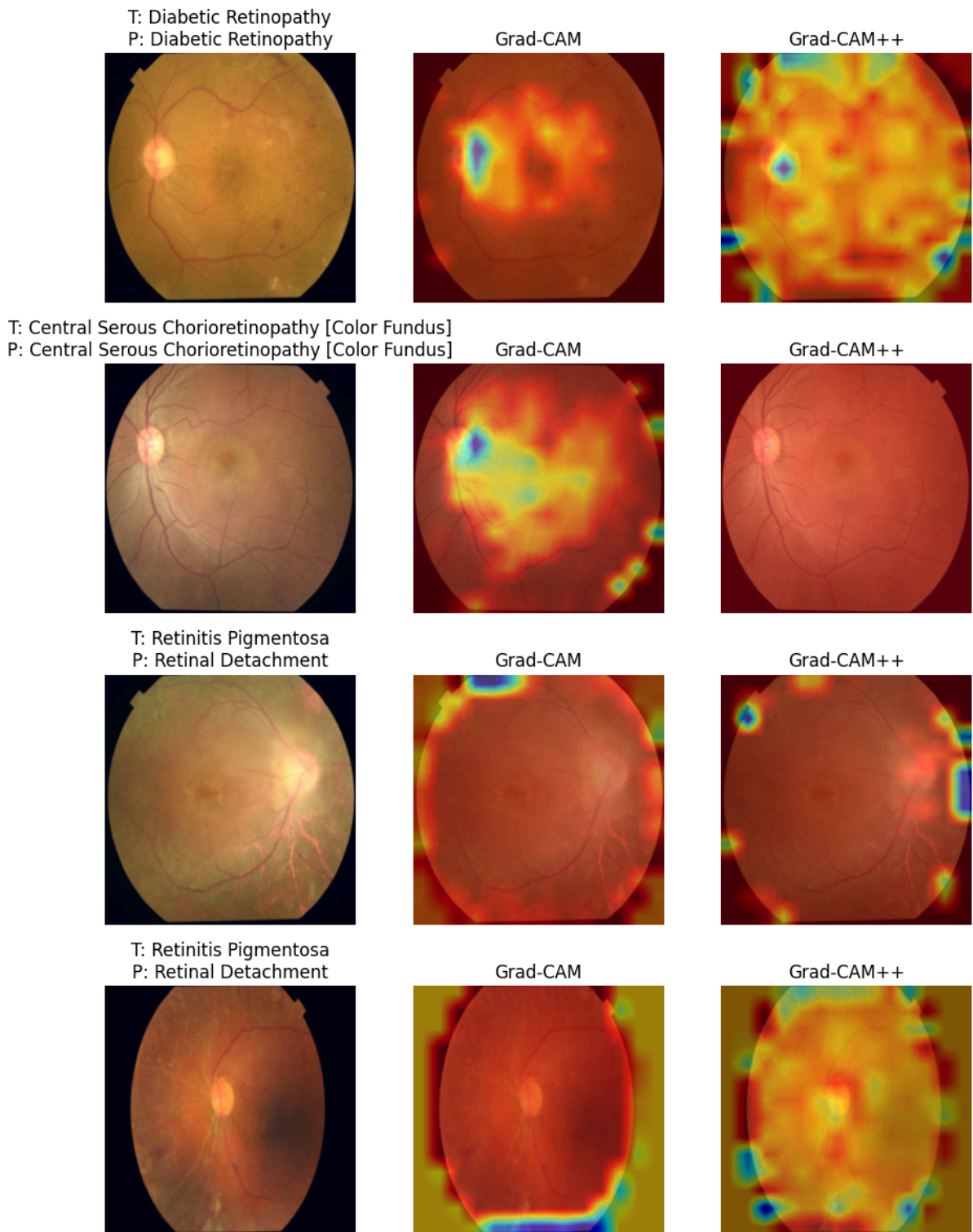
# Interpreting the Unseen: Grad-CAM Visualizations 🔥 👁

To gain deeper insight into how our model makes decisions, we incorporated **Grad-CAM** and **Grad-CAM++** visualizations. These techniques allow us to **highlight the exact regions** in fundus images that influenced the model's predictions the most.

As seen in the figure below, Grad-CAM overlays for different retinal diseases clearly focus on lesion-specific areas — confirming that the model is learning **meaningful clinical features** rather than random patterns.

*For instance,* in cases of **Diabetic Retinopathy**, attention is sharply drawn towards hemorrhages and microaneurysms, while for **Central Serous Chorioretinopathy**, the model learns to focus near the retinal pigment epithelium — closely matching ophthalmologist expectations.

These visualizations enhance model transparency and build confidence in its clinical usability.

T: Diabetic Retinopathy
P: Diabetic Retinopathy

Grad-CAM

Grad-CAM++

T: Central Serous Chorioretinopathy [Color Fundus]
P: Central Serous Chorioretinopathy [Color Fundus]

Grad-CAM

Grad-CAM++

T: Retinitis Pigmentosa
P: Retinal Detachment

Grad-CAM

Grad-CAM++

T: Retinitis Pigmentosa
P: Retinal Detachment

Grad-CAM

Grad-CAM++

# Future Horizons: Multi-Modal Magic ✨📊

While fundus images offer powerful cues, **integrating multiple data modalities** like:

- **OCT scans** (depth view of retinal layers),

- **Patient history** (e.g., diabetes, hypertension),

- **Clinical measurements** (e.g., intraocular pressure)

could significantly enrich the model's decision-making power.

We envision **future versions** of ReSeeAI to adopt a **multi-input transformer architecture** that fuses visual and clinical data at an embedding level — allowing much **earlier detection** of subtle disease patterns not visible through fundus images alone.

## Gradio Demo: Bringing AI to Your Fingertips ⚡🖥️

The trained model is lightweight and designed to be **easily deployed through a simple Gradio web app**.

In future updates, users could upload a fundus image on a browser, and the model would:

- Instantly predict the disease class,
- Show confidence scores for all possible conditions,
- Display a Grad-CAM heatmap highlighting critical regions.

This makes ReSeeAI not just a research prototype, but a **potential real-world screening tool** — accessible even from mobile devices.

# 6. Shaping the Future: Lessons Learned and New Horizons

The journey of developing **ReSeeAI** was filled with learning, experimentation, and perseverance.

Starting with limited knowledge of retinal imaging and large foundation models, we gradually built expertise by diving deep into the capabilities of **RETFound**, one of the first large-scale foundation models trained for ophthalmology.

**Struggles Along the Way:**

- Early models trained on unbalanced datasets achieved modest results (~64% accuracy).
- Several challenges like underperforming classes ("Pterygium", "CSC", "Glaucoma") and augmentation missteps (MixUp/CutMix hurting performance) forced **multiple rounds of trial and error**.
- Managing GPU resources carefully while fine-tuning large ViT backbones was another real-world hurdle.

**Breakthroughs and Key Insights:**

- **Full Fine-Tuning** emerged as the strongest strategy, unlocking ~92% accuracy after careful **class balancing** and **hyperparameter tuning**.
- Dropping classes with too little data drastically stabilized training and improved generalization.
- **Grad-CAM explainability** reassured us that the model focused meaningfully on disease regions, not just random features.

**Where We Head Next:**

- **Multi-modal expansion**: Combine Fundus + OCT + Clinical data for richer, real-world diagnosis.

- **External validation**: Test ReSeeAI on unseen hospital datasets for stronger clinical credibility.

- **Web deployment**: Launch a Gradio-based web app for real-time global accessibility.

- **Fine-grained classification**: Push towards predicting disease stages, not just categories.

---

**ReSeeAI** now stands not just as a technical project —

**but as a small step toward making expert-level eye diagnosis accessible to all.**

🌟 *"Where perseverance meets vision, innovation is born."*

---

# 📚 References and Links Used in Project

## 🔬 Research Papers and Models

- **Retinal eye Detection CNN**

- **RETFound Paper**: Self-supervised retinal foundation model 🖼️ Datasets

- **Fundus Dataset** (color fundus images — created and uploaded by you):

  ➔ https://huggingface.co/datasets/Peacein/color-fundus-eye

- **OCT Dataset** (publicly available retinal OCT images):

  ➔ https://huggingface.co/datasets/MaybeRichard/OCT-retina-classification-2017

---