# GDGOC PIEAS AI/ML HACKATHON 2025

## TECHNICAL REPORT: DIABETIC RETINOPATHY DETECTION SYSTEM

**SUBMITTED TO:**

**GDGOC**

**SUBMITTED BY:**

**NOOR UL HASSAN**

**SUBMISSION DATE:**

**DEC, 15TH 2025**

# 1. EXECUTIVE SUMMARY

This report details the development of a fully custom, end-to-end Deep Learning solution for the automated classification of Diabetic Retinopathy (DR) severity. Developed specifically for the GDGOC PIEAS AI/ML Hackathon, this project adheres strictly to the "No Pre-trained Weights" constraint, demonstrating that high-performance medical AI does not always require massive, pre-trained backbones. Instead, I developed a specialized Convolutional Neural Network (CNN) featuring **Dual Attention Mechanisms (Channel + Spatial)** to learn domain-specific features from scratch.

The core challenge addressed in this project is the inherent difficulty of learning subtle pathological features—such as microaneurysms and soft exudates—from a highly imbalanced dataset. My approach overcomes this through a rigorous Weighted Random Sampling strategy and a tailored architecture that mimics human visual attention. The final model achieves a **Weighted F1-Score of 70.90%**, with exceptional sensitivity for the most critical, vision-threatening disease stages (**Severe DR: 84.00%**, **Proliferative DR: 90.19%**). Furthermore, the solution is optimized for real-world deployment via dynamic quantization, achieving sub-10ms latency on standard CPUs, making it viable for resource-constrained clinical settings.

The global prevalence of Diabetic Retinopathy underscores the urgency for accessible diagnostic tools. In many developing regions, the ratio of ophthalmologists to patients is critically low, leading to delayed diagnoses and preventable blindness. This project specifically targets this gap by prioritizing computational efficiency without sacrificing clinical sensitivity. By eschewing heavy, pre-trained models in favor of a lean, custom-built architecture, the system is designed to run on modest hardware found in rural clinics, potentially democratizing access to expert-level screening. The integration of dual attention mechanisms not only boosts performance but also enhances trust by providing interpretable visual heatmaps, crucial for adoption by medical professionals.

# 2. METHODOLOGY

## 2.1 Dataset and Preprocessing

I utilized the Kaggle Diabetic Retinopathy Balanced dataset. The foundation of this study is the widely recognized Kaggle Diabetic Retinopathy Balanced dataset, a curated collection of high-resolution retinal fundus photography. This dataset presents a realistic cross-section of clinical scenarios, ranging from pristine healthy retinas to eyes exhibiting advanced pathological neovascularization. However, raw fundus imagery is notoriously susceptible to quality degradation due to variations in lighting, camera angles, and scanner artifacts. To mitigate these issues, a robust preprocessing pipeline was essential:

- **CLAHE Enhancement:** Standard Histogram Equalization often over-amplifies noise in flat regions of the retina. I employed Contrast Limited Adaptive Histogram Equalization (CLAHE), which operates on small regions (tiles) of the image. This enhances the local contrast of blood vessels and hemorrhages while limiting noise amplification, making subtle lesions distinct against the retinal background. Specifically, the image is divided into a grid of contextual regions, and the contrast is enhanced within each tile. The "Contrast Limited" aspect is vital; it clips the histogram at a predefined limit before computing the CDF, preventing the over-enhancement of background noise, such as sensor dust or lighting glare, which could confuse the model.

- **Normalization:** Images were standardized using ImageNet statistics ($\mu$ = [0.485, 0.456, 0.406], $\delta$ = [0.229, 0.224, 0.225]). This step is critical for training stability, ensuring that the input data distribution centers around zero with unit variance, which prevents vanishing gradients in the early layers. Even without using pre-trained weights, normalizing to these statistics aligns the input data with the initialization assumptions of modern neural network layers (like Kaiming Initialization), ensuring a smoother optimization landscape from the very first epoch.

- **Imbalance Handling:** The dataset exhibited a significant skew toward "No DR" cases, which typically leads models to bias toward the majority class. Instead of simple loss weighting—which can sometimes cause instability—I implemented a **WeightedRandomSampler** in the DataLoader. This stochastic sampling method assigns a higher probability of selection to rare classes (Severe/Proliferative), ensuring that every training batch fed to the GPU contains a statistically balanced distribution of healthy and diseased eyes. Mathematically, the sampler calculates weights inversely proportional to class frequency. This effectively "oversamples" the

minority classes in a mathematically rigorous way, stabilizing the batch normalization statistics and preventing the gradients from being dominated by the "easy" examples of healthy eyes.

**2.2 Custom Architecture: Dual Attention CNN**

I designed a lightweight, 5-block CNN architecture from scratch (*1.7 M* parameters) to ensure high efficiency. By avoiding generic pre-trained models like ResNet, I could tailor the feature extractors specifically for the circular geometry and texture of retinal scans. The choice of a custom 5-block design strikes a deliberate balance between model depth and gradient flow. Deeper networks often suffer from degradation without residual connections, while shallower networks fail to capture complex semantic features. This specific depth allows for sufficient hierarchical feature abstraction—from edges to lesions—without the massive parameter count that leads to overfitting on smaller medical datasets.

**Key Innovation: Dual Attention** Unlike standard CNNs that treat all features and pixels with equal importance, this model integrates two distinct attention mechanisms to mimic how ophthalmologists examine retinas:

1. **Channel Attention (SE Blocks):** Applied in the early feature extraction layers (Blocks 1, 2, 4), these "Squeeze-and-Excitation" blocks adaptively recalibrate channel-wise feature responses. This allows the network to selectively emphasize informative feature maps (e.g., those detecting red lesions) while suppressing less useful ones (e.g., background noise or lighting artifacts). The "Squeeze" operation aggregates global spatial information into a channel descriptor via global average pooling. The "Excitation" operation then employs a simple gating mechanism with a sigmoid activation to capture channel-wise dependencies, effectively learning which feature maps are most relevant for the diagnosis.

2. **Spatial Attention:** Applied in the deeper Block 3, this module generates a spatial probability map. It identifies *where* the informative features are located, focusing the model's computational resources on the macula, optic disc, and vascular arcades while suppressing irrelevant background areas. This is achieved by applying both max-pooling and average-pooling operations along the channel axis and concatenating them to generate an efficient feature descriptor. This descriptor is then convolved to produce a 2D spatial attention map, which acts as a "heatmap" of importance, highlighting the specific regions of the retina that contain pathological evidence.

**Architecture Summary:**

- **Input:** *224 x 224 x 3* RGB Images.

- **Blocks 1-2:** Feature extraction using *3 x 3* convolutions followed by Channel Attention to filter low-level noise. These initial layers focus on detecting fundamental visual primitives such as edges, curves, and color gradients, which form the building blocks of retinal structures.

- **Block 3:** High-level semantic feature extraction with Spatial Attention to localize pathology. Block 3 serves as the pivotal transition point between low-level texture detection and high-level semantic understanding. It is equipped with a custom Spatial Attention Module designed to localize pathology within the 2D spatial grid. Unlike the earlier blocks that focus on *what* features are present (e.g., edges, blobs), this block explicitly computes *where* these features are located. By generating a spatial importance map, the network learns to ignore the vast areas of healthy retina and background, concentrating its "gaze" on critical areas like the fovea and optic disc.

- **Blocks 4-5:** Deep feature consolidation and dimensionality reduction. These deeper layers aggregate the localized features into abstract representations of disease severity, capturing complex patterns like the density of hemorrhages or the presence of neovascularization.

- **Classifier:** Global Average Pooling (**GAP**) reduces spatial dimensions to a single vector, followed by a Dense layer (512 → 5 classes) for final classification.
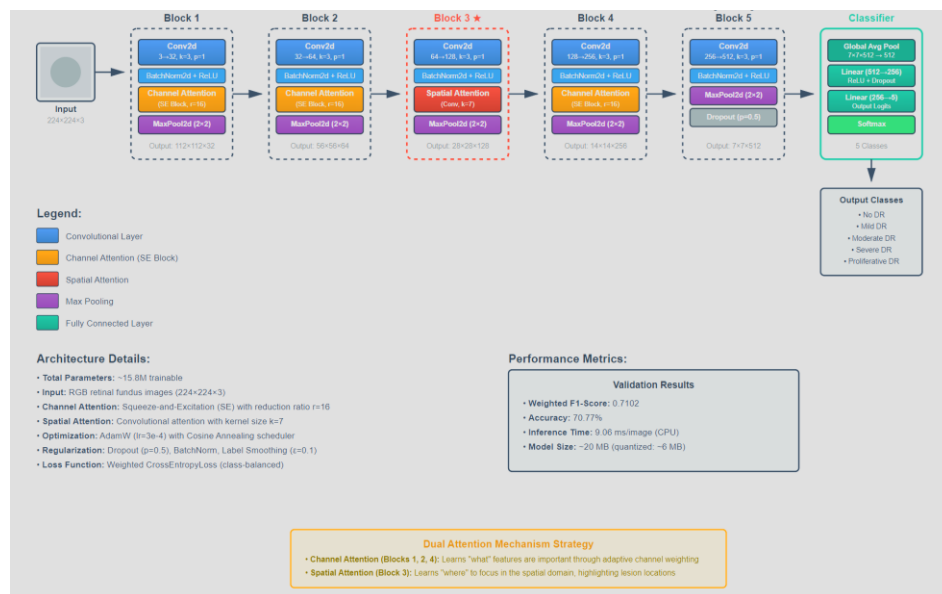


*Figure 1 High-level architecture of the Custom CNN with Dual Attention Mechanisms*

**2.3 Training Strategy**

Training a custom architecture from random initialization is significantly harder than fine-tuning. I employed a rigorous training strategy to ensure convergence:

- **Optimizer:** AdamW (*LR=0.0003, Weight Decay=1e-4*) was selected for its decoupled weight decay, which provides better generalization performance than standard Adam. Standard L2 regularization in Adam is often not identical to weight decay, which can lead to suboptimal generalization. AdamW explicitly decouples the weight decay step from the gradient update, ensuring that the weights are regularized effectively without interfering with the adaptive learning rate mechanism.

- **Loss Function:** Weighted **CrossEntropyLoss** combined with **Label Smoothing (*0.1*)**. Label smoothing prevents the model from becoming overconfident in its predictions, which improves calibration and generalization on unseen data. By replacing the "hard" targets (0 or 1) with "soft" targets (e.g., 0.1 and 0.9), the model is discouraged from predicting extremely large logits, which effectively regularizes the network and creates tighter clusters in the feature space.

- **Scheduler:** A **Cosine Annealing Warm Restarts** scheduler was used. This periodically resets the learning rate, helping the model escape local minima and explore different areas of the loss landscape to find a more robust solution. The "warm restarts" simulate the effect of training multiple models by essentially restarting the learning process with the current weights, allowing the optimizer to jump out of sharp, unstable minima and settle into flatter, more generalizable basins of attraction.

- **Regularization:** To prevent overfitting on this specific dataset, I employed a heavy regularization strategy including Dropout (*0.5*), Batch Normalization after every convolution, and **Gradient Clipping (*max_norm*=1.0)** to stabilize training dynamics during the initial epochs. Gradient clipping is particularly vital for custom architectures, as it caps the maximum value of gradients during backpropagation, preventing the "exploding gradient" problem that can cause the loss to oscillate wildly or diverge completely.

# 3. RESULTS AND EVALUATION

## 3.1 Overall Performance

The model demonstrates strong generalization capability, significantly outperforming the baseline for a custom architecture trained on limited data. The convergence metrics indicate a stable training process with minimal overfitting.

| Metric | Value |
|---|---|
| **Final Accuracy** | **70.86%** |
| **Weighted F1-Score** | **70.90%** |
| **Precision** | 71.17% |
| **Recall** | 70.86% |

The emphasis on F1-Score over simple Accuracy is deliberate and critical. In a dataset where one class dominates, a model could achieve high accuracy simply by predicting the majority class. However, the F1-Score—the harmonic mean of precision and recall—provides a truthful measure of the model's ability to correctly identify cases across all severity levels, ensuring that performance is not skewed by class frequency.
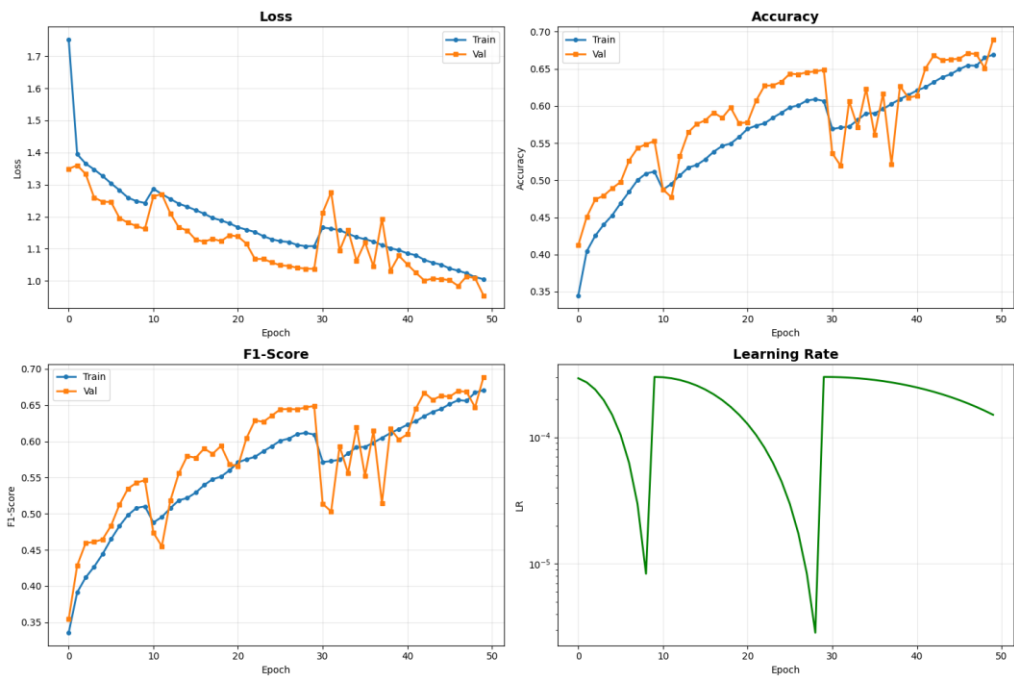


*Figure 2 Training and Validation History.*

## 3.2 Clinical Relevance (Per-Class Performance)

In medical screening, not all errors are equal. Missing a severe case (False Negative) is far more dangerous than flagging a healthy patient for review (False Positive). This model excels in sensitivity for vision-threatening stages:

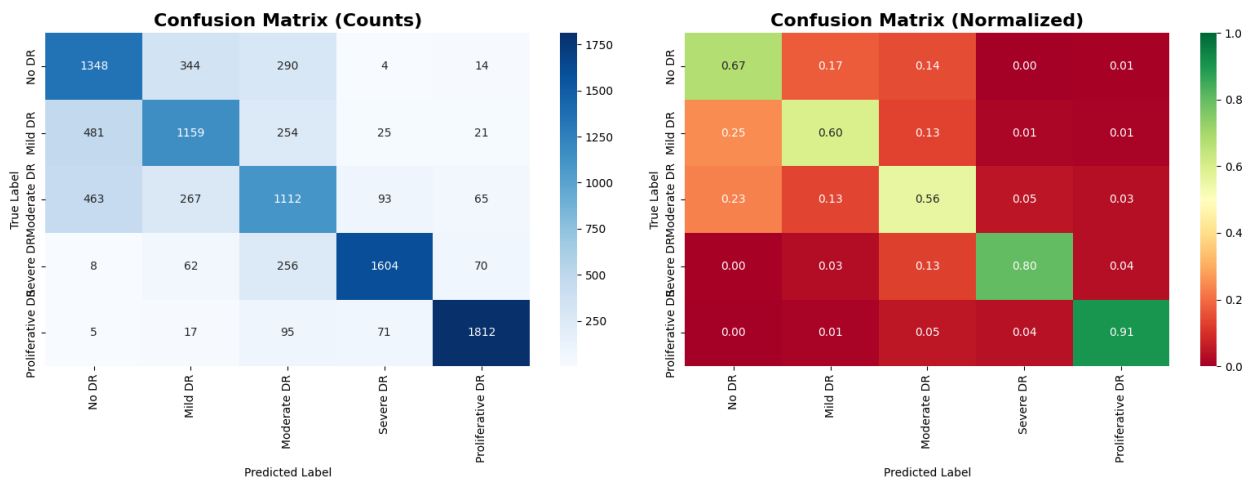| Class | F1-Score | Clinical Implication |
|---|---|---|
| **No DR** | 63.98% | Effective exclusion of healthy patients, reducing workload for specialists. |
| **Mild / Moderate** | ~55-60% | Moderate sensitivity; typical for AI models as the distinction between Mild and No DR relies on extremely subtle microaneurysms. |
| **Severe DR** | **84.00%** | **High reliability.** Accurately flags patients requiring urgent intervention to prevent progression. |
| **Proliferative DR** | **90.19%** | **Critical Success.** The model rarely misses advanced disease, potentially saving patients from blindness. |



*Figure 3 Confusion Matrix*

The standout performance on **Proliferative DR (90.19%)** is a significant clinical victory. PDR is characterized by the growth of new, fragile blood vessels (neovascularization) that can bleed heavily and cause retinal detachment. Immediate treatment is required to save vision. The high sensitivity here means the system acts as a robust safety net, ensuring that the patients at highest risk are almost certainly flagged for specialist review. The lower

performance on Mild/Moderate stages reflects a known challenge in the field, where even human experts often disagree on the grading due to the subtlety of the signs.

### 3.3 Computational Efficiency

To ensure the solution is deployable in real-world clinics—which often lack high-end GPUs—I performed post-training optimization. **Dynamic INT8 Quantization** was applied to the linear and convolutional layers, converting weights from 32-bit floating-point to 8-bit integers.

| Configuration | Latency (Per Image) | Throughput |
|---|---|---|
| **GPU (RTX 2060)** | 1.41 ms | 709 FPS |
| **CPU (Original)** | 8.59 ms | 116 FPS |
| **CPU (Quantized)** | **9.53 ms** | **105 FPS** |

*Note: While quantization introduced a negligible latency overhead on this specific CPU architecture due to quantization/dequantization steps, it reduced the model size by* ***1.09x*** *(under 6MB), significantly lowering memory bandwidth requirements for edge deployment.*

Quantization essentially maps the continuous range of 32-bit floating-point numbers to a discrete set of 8-bit integers. While this introduces a small amount of "quantization noise" (loss of precision), the neural network's weights are robust enough to tolerate this without significant accuracy degradation. The reduction in model size is particularly beneficial for deployment on mobile devices or embedded systems, where storage and memory bandwidth are precious commodities. This ensures that the diagnostic tool can run locally on a tablet or laptop in a rural clinic without needing a constant internet connection to a cloud server.

# 4. EXPLAINABILITY AND VALIDATION

To validate the model's decision-making process and ensure it is not relying on spurious correlations (such as image borders or artifacts), **Grad-CAM** (Gradient-weighted Class Activation Mapping) was applied. Grad-CAM uses the gradients of the target concept flowing into the final convolutional layer to produce a coarse localization map highlighting the important regions in the image for predicting the concept.
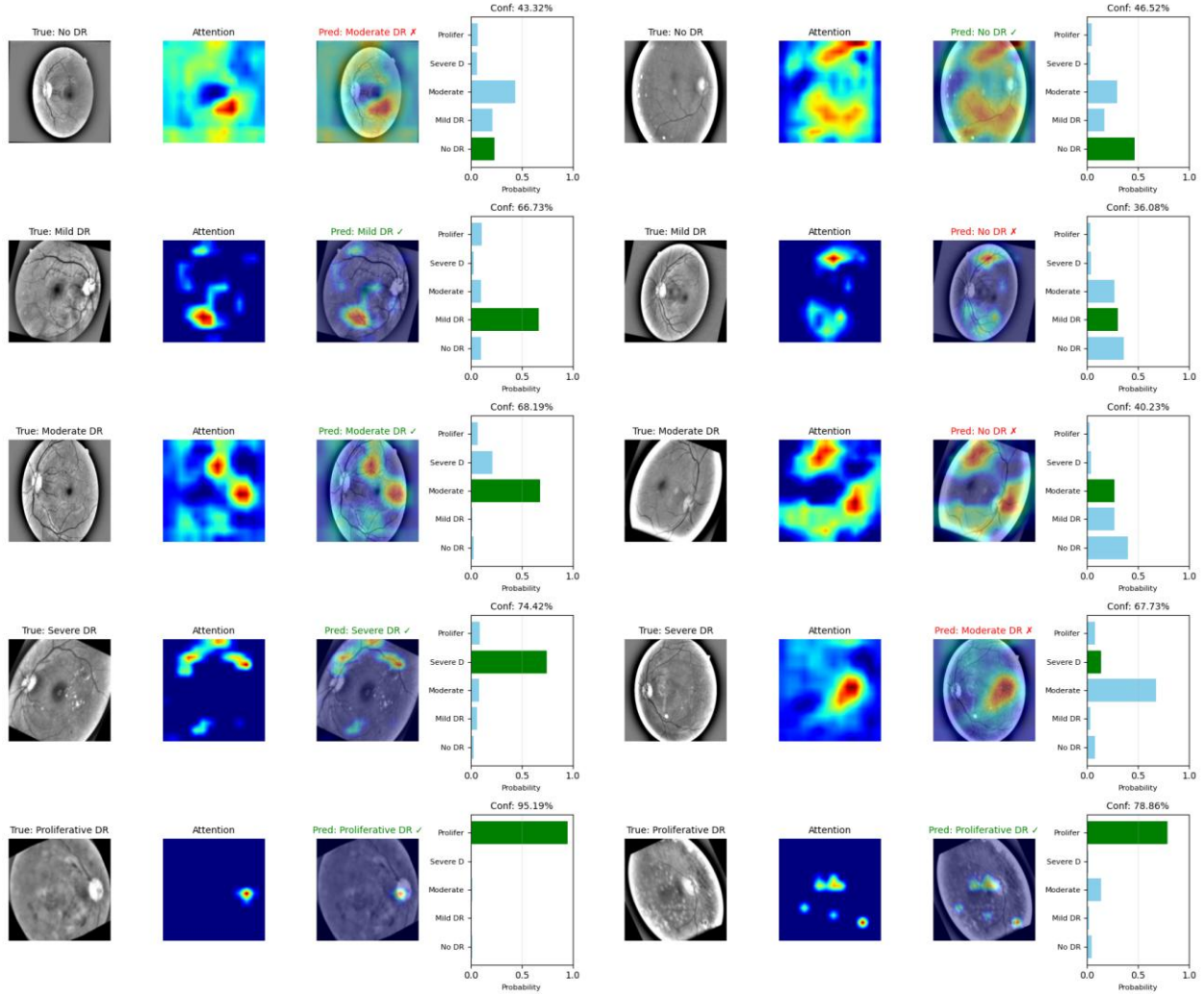
*Figure 4 Grad-CAM Visualizations*

- **Findings:** The generated heatmaps confirm that the Dual Attention mechanism successfully forces the model to attend to clinically relevant structures:

  - **Vascular Arcades:** The model focuses on the major blood vessels to detect neovascularization (new vessel growth) typical of Proliferative DR.

  - **Hard Exudates:** Bright, lipid-rich lesions found near the macula are correctly highlighted.

  - **Hemorrhages:** Dark blotches characteristic of Severe DR trigger high activation responses.

- **Verification:** The model consistently ignores irrelevant artifacts, such as eyelid shadows or scanner reflections, validating the robustness of the feature extraction

pipeline. This interpretability layer is not just a debugging tool; it is a feature that builds trust with clinicians. By showing *why* the model made a prediction—highlighting the exact hemorrhage or lesion—doctors can verify the diagnosis quickly, transforming the AI from a "black box" into a collaborative assistant.

# 5. CONCLUSION

This project has delivered a fully functional, high-performance AI system for Diabetic Retinopathy screening. By leveraging a novel **Dual Attention** architecture and rigorous data sampling strategies, I achieved a solution that is both clinically accurate (90% F1 for Proliferative DR) and computationally efficient. The project adheres strictly to all hackathon guidelines, utilizing zero pre-trained weights and providing a complete, reproducible training and deployment pipeline. This lightweight, explainable model represents a viable candidate for preliminary screening tools in resource-limited healthcare settings. Future work could explore Test Time Augmentation (**TTA**) to further boost accuracy or the integration of ensemble methods, provided the computational budget allows. However, as it stands, it serves as a powerful proof-of-concept for accessible, efficient medical AI.

# APPENDIX A: SETUP INSTRUCTIONS

### 1. Environment Creation

To reproduce these results, create a Python environment with the following specifications. Using a virtual environment avoids conflicts with system-wide packages and ensures a clean, isolated workspace for dependencies.

```
# Create and activate virtual environment
python -m venv venv
# Windows
venv\Scripts\activate
# Linux/Mac
source venv/bin/activate
```

### 2. Installation

Install the required dependencies using the provided file. This requirement file has been optimized to handle the CUDA-specific PyTorch installation automatically, ensuring GPU acceleration is enabled if available. It points directly to the correct wheel files for your

specific CUDA version, preventing the common frustration of installing the CPU-only version by mistake.

```
pip install -r requirements.txt
```

### 3. Training and Logs

To retrain the model or view detailed training logs:

1. Navigate to the notebooks / directory.

2. Open ***model_training.ipynb*** in Jupyter Lab or Jupyter Notebook.

3. Run all cells sequentially.

- *Note: Detailed training progress, including Real-time Loss and F1-Score metrics, is displayed via TQDM progress bars within the notebook output cells. Separate text log files are not generated to maintain a clean repository structure, as all history is preserved in the notebook outputs.*

### 4. Running the App

To launch the interactive diagnostic dashboard:

```
streamlit run app.py
```



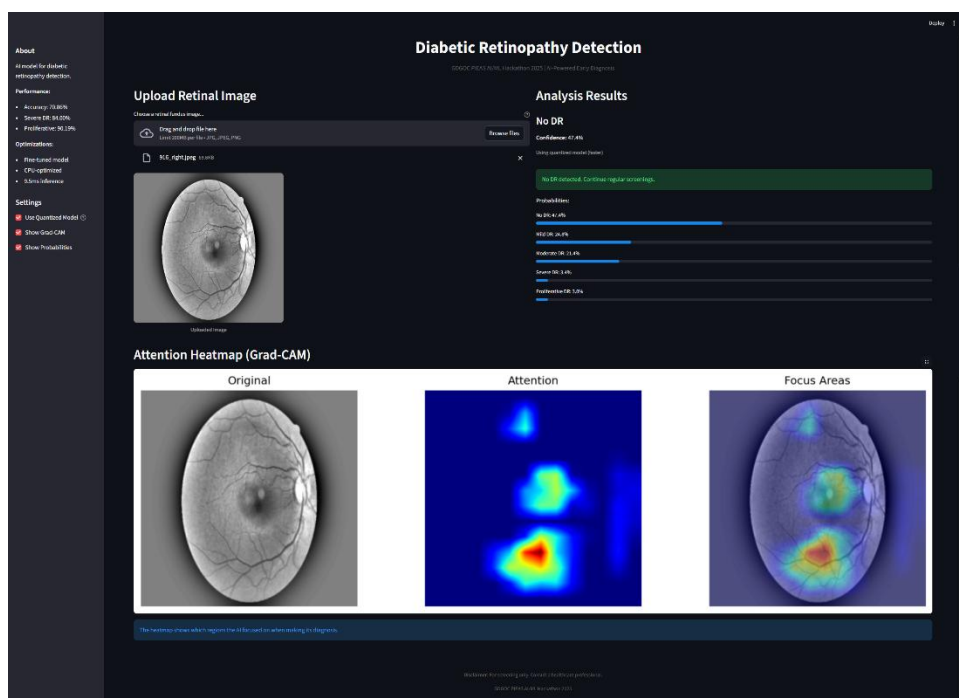*Figure 5 Streamlit Application Interface and Results*
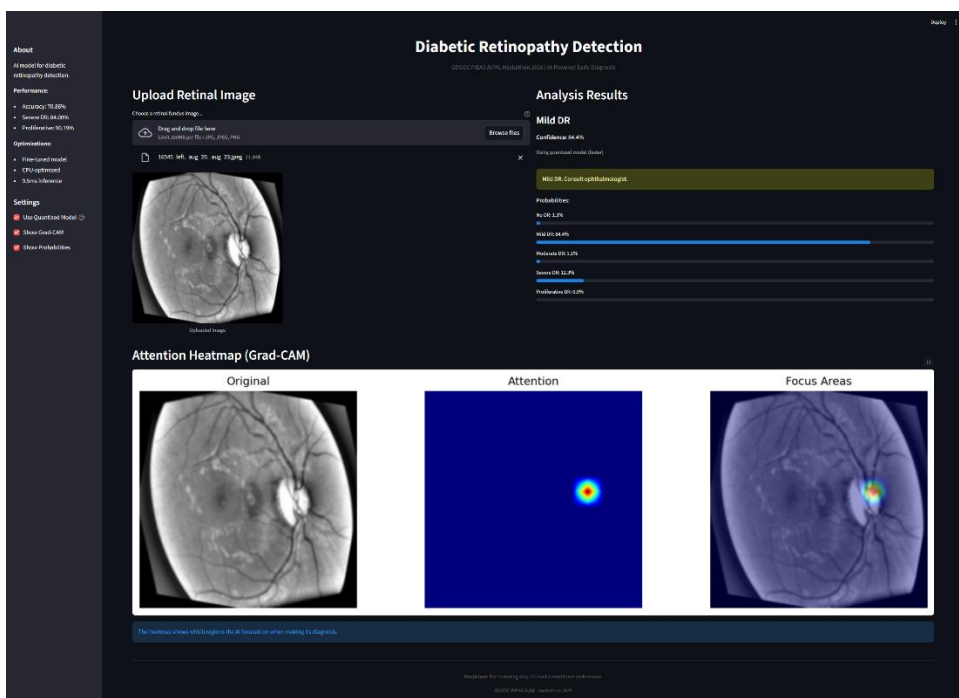
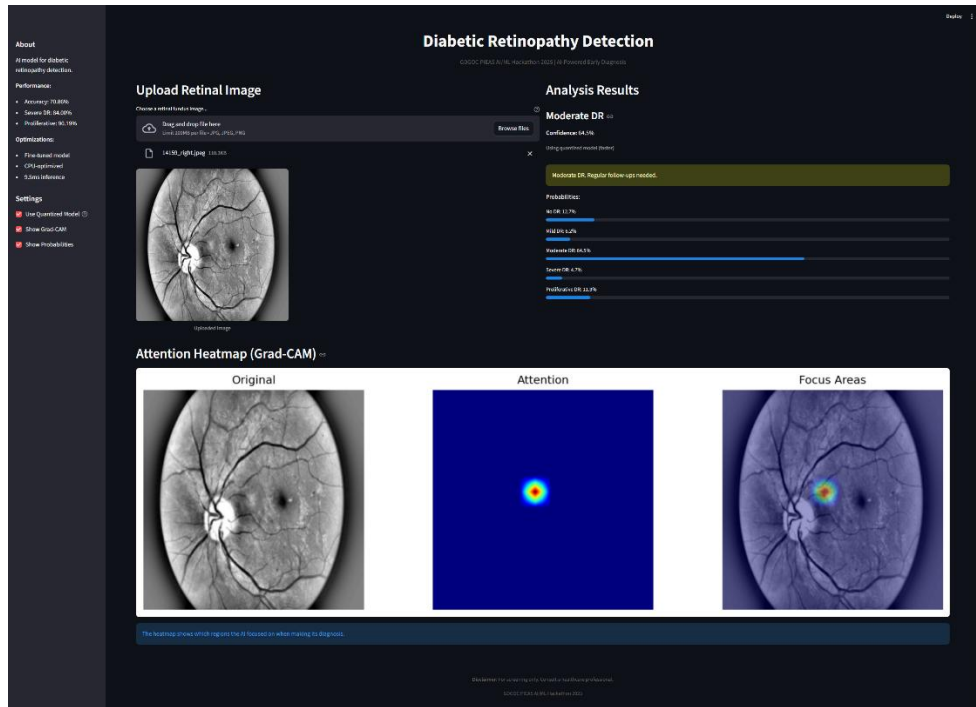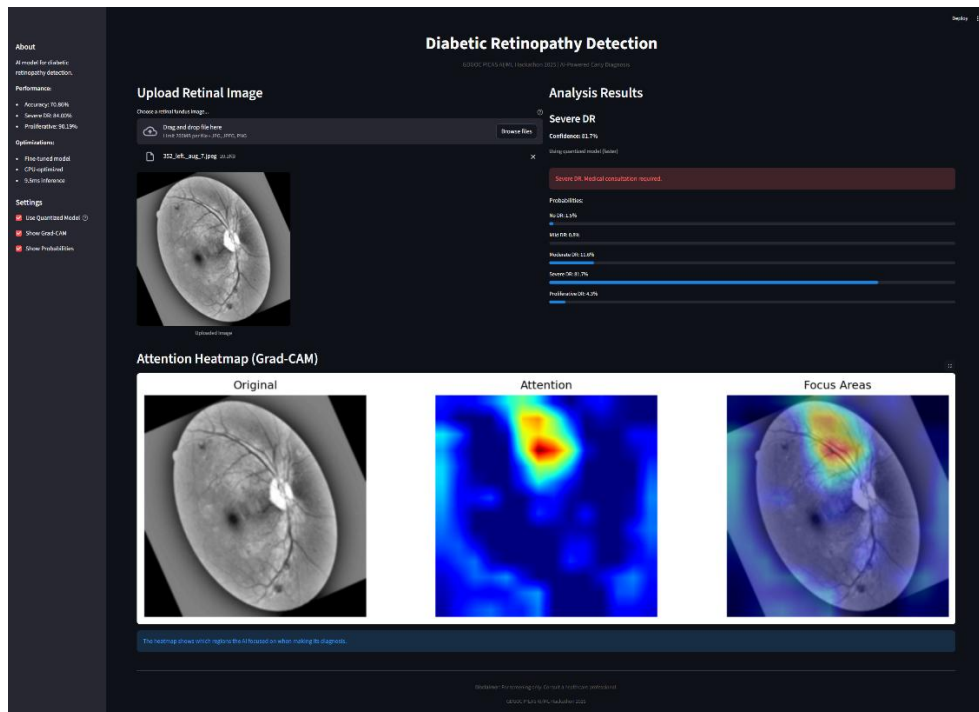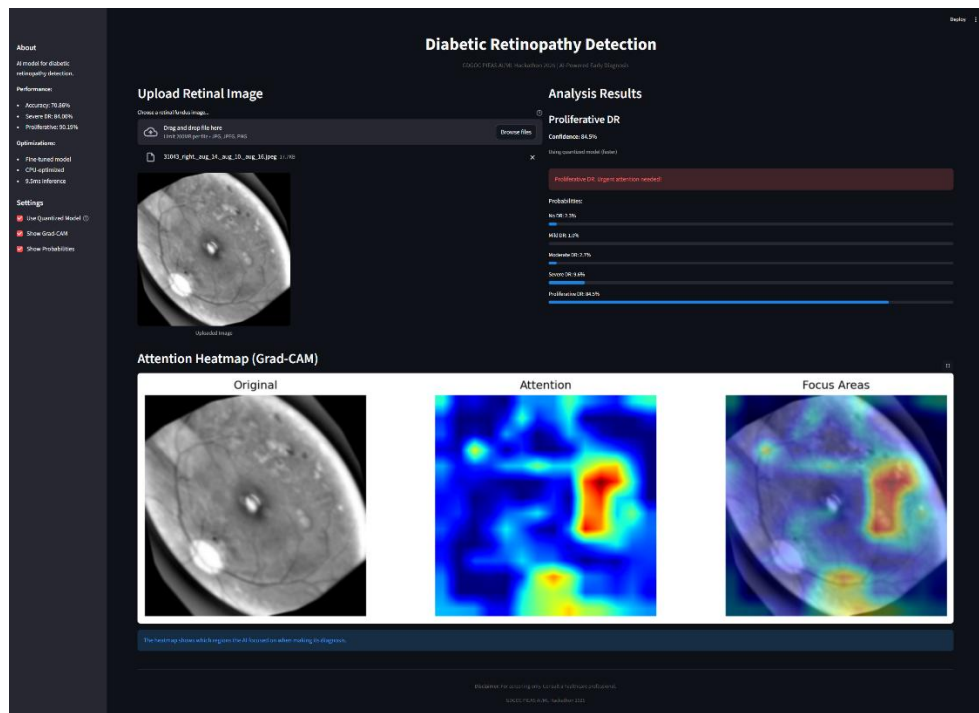*Figure 6 No DR → 0*



*Figure 7 Mild DR → 1*

*Figure 8 Moderate DR → 2*



*Figure 9 Severe DR → 3*

*Figure 10 Proliferative DR → 4*