

Automated Grading of Diabetic Retinopathy Using a Multi-Scale Deep Convolutional Neural Network with Attention Mechanism

KM Sonam¹, Prof. Jagdeep Singh²

¹M.Tech, Computer Science and Engineering
Shree Vanktshvera University Gajrola, Uttar Pradesh, India

²Supervisor, Department of Computer Science and Engineering, Shree Vanktshvera University Gajrola

ABSTRACT

Diabetic Retinopathy (DR) is a leading cause of preventable blindness, affecting over 93 million individuals worldwide. Early-stage detection remains critically challenging due to subtle fundus lesion morphologies and the scarcity of trained ophthalmologists in low-resource settings. This paper proposes a novel Multi-Scale Deep Convolutional Neural Network augmented with a Channel-Spatial Attention Mechanism (MS-CNN-CSAM) for automated, end-to-end DR severity grading. The proposed architecture integrates multi-scale feature extraction through parallel convolutional streams at three distinct receptive field sizes (3×3, 5×5, 9×9), a dual-branch attention module based on Squeeze-and-Excitation (SE) and Convolutional Block Attention Module (CBAM) principles, and a classification head employing Focal Loss to handle severe class imbalance inherent in clinical DR datasets. Experiments are conducted on two publicly available benchmark datasets: the Kaggle EyePACS Dataset (88,702 fundus images) and the APTOS 2019 Blindness Detection dataset (3,662 images). The proposed model achieves a five-class grading accuracy of 93.7% (Quadratic Weighted Kappa: 0.934) on EyePACS and 95.2% (QWK: 0.947) on APTOS 2019, surpassing existing state-of-the-art approaches including VGG-19 (88.4%), ResNet-50 (90.1%), Inception-V3 (91.2%), and EfficientNet-B4 (92.8%). A comprehensive ablation study confirms that the attention mechanism and multi-scale fusion contribute 2.3% and 1.8% incremental accuracy gains, respectively. The proposed model demonstrates strong generalizability, computational efficiency suitable for GPU-constrained clinical environments, and statistical significance under McNemar's test ($p < 0.001$). The work establishes a clinically viable AI-driven screening pipeline that can meaningfully reduce the global burden of preventable DR-induced vision loss.

Keywords: Diabetic Retinopathy, Convolutional Neural Network, Attention Mechanism, Multi-Scale Feature Extraction, Fundus Image Analysis, Deep Learning, Medical Image Classification, Class Imbalance, Quadratic Weighted Kappa

1. INTRODUCTION

1.1 Background of Artificial Intelligence and Machine Learning

Artificial Intelligence (AI) has undergone a transformative renaissance over the past decade, driven by exponential increases in computational power, the availability of large-scale annotated datasets, and fundamental algorithmic breakthroughs. Within the AI paradigm, Machine Learning (ML) constitutes the computational framework wherein systems learn statistical patterns directly from empirical data without explicit rule-based programming [1]. The seminal work of LeCun et al. [2] on Gradient-Based Learning Applied to Document Recognition established the theoretical foundations of convolutional architectures that now underpin virtually all state-of-the-art visual recognition systems. Deep Learning, a sub-discipline of ML characterized by hierarchical feature representations learned through multi-layered neural network architectures, has demonstrated superhuman performance in image classification [3], natural language processing [4], and speech recognition [5].

In the biomedical domain specifically, AI and ML methodologies have catalyzed paradigmatic shifts in diagnostic imaging, drug discovery, genomics, and clinical decision support. The integration of convolutional neural networks into ophthalmological screening workflows represents one of the most compelling demonstrations of AI's potential to democratize access to specialized medical expertise.

1.2 Neural Networks in Modern Applications

Artificial Neural Networks (ANNs) are computational models abstractly inspired by biological neural architectures in the mammalian cortex. Their ability to approximate arbitrarily complex non-linear functions through the composition of differentiable transformations renders them uniquely suited for high-dimensional perceptual tasks including image segmentation, object detection, and medical grade classification [6]. Contemporary deep learning frameworks such as TensorFlow [7] and PyTorch [8] have substantially reduced the engineering barrier to training large-scale neural networks, enabling widespread application in clinical and research settings. Specific architectures—including Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), Transformers, and hybrid models—have each demonstrated domain-specific advantages that motivate the exploration of tailored architectural designs for ophthalmological image analysis.

1.3 Problem Statement

Diabetic Retinopathy (DR) is a microvascular complication of diabetes mellitus that progressively damages the retinal vasculature, ultimately culminating in complete vision loss if untreated. The International Diabetes Federation estimates that 537 million adults currently live with diabetes globally, with projections reaching 783 million by 2045 [9]. Of these, approximately 35% will develop some degree of DR. The severity of DR is clinically stratified on a five-level International Clinical Diabetic Retinopathy (ICDR) severity scale: No DR (Grade 0), Mild NPDR (Grade 1), Moderate NPDR (Grade 2), Severe NPDR (Grade 3), and Proliferative DR (Grade 4).

Manual grading by trained retinal specialists is the gold standard for DR diagnosis. However, this approach suffers from three critical operational constraints: (i) the severe global shortage of ophthalmologists, particularly in South Asian and Sub-Saharan African regions with the highest diabetes prevalence; (ii) significant inter-rater and intra-rater grading variability (reported kappa values of 0.64–0.81 in clinical

studies); and (iii) the infeasibility of screening all at-risk individuals at the recommended annual frequency given workforce constraints. Automated DR grading systems therefore represent an urgent clinical imperative.

1.4 Motivation

My interest in this problem grew out of coursework in Machine Learning and Pattern Recognition, where I kept returning to a question that felt unresolved: why do models that perform well on standard benchmarks still struggle to be trusted in real clinical environments? As I surveyed the DR detection literature, the answer became increasingly clear — not because any single paper was technically weak, but because each paper solved one piece of the puzzle in isolation while leaving the others untouched. One paper would improve accuracy overall but ignore the fact that 73% of its training data belonged to a single class. Another would introduce an attention mechanism but bolt it onto a backbone designed for ImageNet, not for retinal images. A third would report impressive AUC figures but on a dataset so clean it bore little resemblance to images collected at a rural screening camp. What none of them did was treat the problem as a system — where the data imbalance strategy, the feature extraction design, and the attention mechanism all needed to be designed together to work together. That gap is what motivated this research. The goal was not to invent a fundamentally new algorithm, but to ask a different question: what happens if you combine the best available techniques — multi-scale convolution, dual-branch attention, and imbalance-aware training — into one coherent, end-to-end pipeline that is honest about each component's contribution?

1.5 Research Objectives

This research is guided by the following primary objectives:

- Design and implement a novel Multi-Scale CNN architecture with parallel convolutional streams that capture micro- and macro-scale retinal features simultaneously.
- Integrate a dual-branch Channel-Spatial Attention Mechanism to focus network attention on clinically significant lesion regions.
- Develop a class-imbalance-robust training strategy using Focal Loss and stratified oversampling to improve grading performance on minority DR severity classes.
- Benchmark the proposed model against established state-of-the-art baselines (VGG-19, ResNet-50, Inception-V3, EfficientNet-B4) on two public datasets.
- Conduct a rigorous ablation study to quantify the individual contribution of each architectural component.
- Evaluate clinical deployment feasibility through inference time and hardware resource profiling.

1.6 Research Contributions

The principal contributions of this work are enumerated as follows:

- [C1] A novel Multi-Scale CNN architecture (MS-CNN-CSAM) specifically designed for multi-scale retinal lesion feature extraction with three parallel convolutional streams.
- [C2] A dual-branch attention module integrating channel attention (SE-Net) and spatial attention (CBAM) for lesion-localized feature recalibration.

- [C3] An empirically validated class-imbalance mitigation strategy combining Focal Loss ($\gamma=2$, $\alpha=0.25$) with Synthetic Minority Oversampling Technique for Fundus Images (SMOTE-FI).
- [C4] A comprehensive performance evaluation achieving state-of-the-art results of 93.7% accuracy and QWK = 0.934 on EyePACS and 95.2% accuracy with QWK = 0.947 on APTOS 2019.
- [C5] Publicly released pre-trained model weights, training configuration files, and evaluation scripts to facilitate reproducibility.

1.7 Organization of the Paper

The remainder of this paper is organized as follows: Section 2 presents a comprehensive literature review encompassing neural network fundamentals, architectural variants, mathematical foundations, and research gap analysis. Section 3 details the proposed MS-CNN-CSAM methodology, dataset description, feature engineering, and training strategy. Section 4 describes the experimental setup. Section 5 presents quantitative results, comparative analysis, and discussion. Sections 6 and 7 discuss the model's advantages and limitations, respectively. Section 8 concludes the paper, and references along with appendices follow thereafter.

2. LITERATURE REVIEW

2.1 Overview of Neural Networks

2.1.1 Definition

An Artificial Neural Network (ANN) is a parameterized, directed computational graph comprising interconnected processing nodes (neurons) organized in layers. Each neuron computes a weighted linear combination of its inputs followed by a non-linear activation function, yielding a transformation that, when composed across multiple layers, can represent highly complex multivariate functions. Formally, a neural network with L layers computes: $f(x) = f_L \circ f_{L-1} \circ \dots \circ f_1(x)$, where each layer f_l applies a learned affine transformation followed by an element-wise nonlinearity [14].

2.2 Types of Neural Networks

2.2.1 Feedforward Neural Network (FNN)

The Feedforward Neural Network, also termed a Multilayer Perceptron (MLP), constitutes the canonical deep learning architecture wherein information propagates unidirectionally from input to output through one or more hidden layers. Each layer performs a non-linear affine transformation: $z^{(l)} = W^{(l)}a^{(l-1)} + b^{(l)}$; $a^{(l)} = \sigma(z^{(l)})$, where $W^{(l)}$ and $b^{(l)}$ denote the weight matrix and bias vector of layer l , and σ is the activation function [14]. FNNs are universal approximators but scale poorly with spatial input dimensionality, motivating convolutional architectures for image data.

2.2.2 Convolutional Neural Network (CNN)

CNNs exploit three key inductive biases—local connectivity, parameter sharing, and translation equivariance—to efficiently process grid-structured data. A convolutional layer applies a bank of learnable filters K of spatial size $k \times k$ to the input feature map F , producing output activations: $O_{\{i,j,m\}} = \sum_u \sum_v \sum_c K_{\{u,v,c,m\}} \cdot F_{\{i+u,j+v,c\}} + b_m$. Landmark architectures include LeNet-5 [2], AlexNet [3],

VGGNet [10], GoogLeNet/Inception [12], ResNet [11], DenseNet [19], MobileNet [20], and EfficientNet [13].

2.2.3 Recurrent Neural Network (RNN)

Recurrent Neural Networks process sequential data through a hidden state h_t that encodes temporal context: $h_t = \sigma(W_h \cdot h_{t-1} + W_x \cdot x_t + b)$. Despite their expressive power, vanilla RNNs suffer from vanishing and exploding gradients during backpropagation through time (BPTT), limiting effective sequence length [21].

2.2.4 Long Short-Term Memory (LSTM)

The LSTM architecture, introduced by Hochreiter and Schmidhuber [22] in 1997, addresses the vanishing gradient problem through a gated cell state mechanism comprising input gate i_t , forget gate f_t , output gate o_t , and cell state c_t : $f_t = \sigma(W_f[h_{t-1}, x_t] + b_f)$; $i_t = \sigma(W_i[h_{t-1}, x_t] + b_i)$; $c_t = f_t \odot c_{t-1} + i_t \odot \tanh(W_c[h_{t-1}, x_t] + b_c)$; $o_t = \sigma(W_o[h_{t-1}, x_t] + b_o)$; $h_t = o_t \odot \tanh(c_t)$. LSTMs are widely applied in time-series medical signal processing, clinical note generation, and longitudinal disease progression modeling.

2.2.5 Autoencoders

Autoencoders are unsupervised neural architectures comprising an encoder $E: x \rightarrow z$ and decoder $D: z \rightarrow \hat{x}$ trained to minimize a reconstruction loss $L(x, \hat{x})$. The latent bottleneck representation z serves as a compact data manifold embedding. Variational Autoencoders (VAEs) [23] extend this to a probabilistic generative framework by regularizing z under a Gaussian prior. In medical imaging, autoencoders are applied for anomaly detection, image denoising, and semi-supervised representation learning.

2.2.6 Generative Adversarial Networks (GANs)

GANs [24] comprise two networks trained adversarially: a generator $G(z)$ that synthesizes realistic samples from random noise z , and a discriminator $D(x)$ that distinguishes real from synthetic samples. The training objective is: $\min_G \max_D E_{x \sim p_{\text{data}}}[\log D(x)] + E_{z \sim p_z}[\log(1 - D(G(z)))]$. In the context of DR, GANs are employed for augmenting minority class training samples, synthesizing pathological fundus images, and enabling domain adaptation across heterogeneous imaging hardware [25].

2.2.7 Transformer-Based Neural Networks

The Transformer architecture [26], originally proposed for natural language processing, employs a multi-head self-attention mechanism that models long-range dependencies: $\text{Attention}(Q, K, V) = \text{softmax}(QK^T / \sqrt{d_k})V$. The Vision Transformer (ViT) [27] adapts this mechanism to image patches, treating each 16×16 patch as a token. Recent hybrid architectures such as Swin Transformer [28] and BEiT [29] have demonstrated competitive performance on medical image benchmarks, motivating their use as backbone components in fundus analysis pipelines.

2.3 Key Components of Neural Networks

2.3.1 Architecture Components

The input layer receives raw feature vectors or pixel-valued tensors without transformation. Hidden layers apply successive learned transformations, with depth encoding representational hierarchy—lower layers encode low-level features (edges, textures) and higher layers encode semantic concepts (lesion

morphology, anatomical structures). The output layer produces class probabilities via softmax: $p(y=k|x) = \exp(z_k) / \sum_j \exp(z_j)$, enabling categorical DR grade prediction.

Weights $W^{(l)} \in \mathbb{R}^{n_l \times n_{l-1}}$ encode the learned linear projection strength; biases $b^{(l)} \in \mathbb{R}^{n_l}$ shift the activation hyperplane. Activation functions introduce non-linearity indispensable for universal function approximation. The Rectified Linear Unit (ReLU): $f(x) = \max(0, x)$, Leaky ReLU: $f(x) = \max(\alpha x, x)$ with $\alpha=0.01$, GELU: $f(x) = x \cdot \Phi(x)$, and Swish: $f(x) = x \cdot \sigma(x)$ are dominant activations in modern architectures.

Loss functions quantify prediction error. Cross-Entropy Loss: $L_{CE} = -\sum_k y_k \log(p_k)$ is standard for multi-class classification, while Focal Loss: $L_{FL} = -\alpha_t (1-p_t)^\gamma \log(p_t)$ down-weights easy negatives to address class imbalance [30]. Optimization algorithms iteratively update parameters to minimize loss: SGD with momentum, RMSProp, and Adam (Adaptive Moment Estimation) are most widely adopted. Adam combines adaptive learning rates with first and second moment estimates: $m_t = \beta_1 m_{t-1} + (1-\beta_1)g_t$; $v_t = \beta_2 v_{t-1} + (1-\beta_2)g_t^2$; $\theta_{t+1} = \theta_t - \eta \cdot \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon)$.

2.4 Mathematical Foundations

2.4.1 Forward Propagation

For a network with L layers, forward propagation computes pre-activation and post-activation tensors layer by layer:

$$z^{(l)} = W^{(l)} \cdot a^{(l-1)} + b^{(l)} \quad [\text{Pre-activation, Layer } l] \quad \dots(1)$$

$$a^{(l)} = \sigma(z^{(l)}) \quad [\text{Post-activation, Layer } l] \quad \dots(2)$$

$$\hat{y} = \text{softmax}(z^{(L)}) = \exp(z^{(L)}) / \|\exp(z^{(L)})\|_1 \quad [\text{Output Probability, Layer } L] \quad \dots(3)$$

2.4.2 Backpropagation

Backpropagation computes the gradient of the scalar loss L with respect to each parameter via the chain rule. The output layer delta is: $\delta^{(L)} = \nabla_{\{a^{(L)}\}L} \odot \sigma'(z^{(L)})$. For hidden layer l : $\delta^{(l)} = (W^{(l+1)})^T \cdot \delta^{(l+1)} \odot \sigma'(z^{(l)})$. Parameter gradients are: $\partial L / \partial W^{(l)} = \delta^{(l)} \cdot (a^{(l-1)})^T$; $\partial L / \partial b^{(l)} = \delta^{(l)}$.

$$\delta^{(l)} = (W^{(l+1)})^T \cdot \delta^{(l+1)} \odot \sigma'(z^{(l)}) \quad \dots(4)$$

$$\partial L / \partial W^{(l)} = \delta^{(l)} \cdot (a^{(l-1)})^T \quad \dots(5)$$

2.4.3 Gradient Descent Optimization

The vanilla gradient descent update rule: $\theta \leftarrow \theta - \eta \nabla_{\theta} L(\theta)$ suffers from high variance in stochastic settings. Adam [31] ameliorates this through bias-corrected adaptive moment estimation:

$$\hat{m}_t = m_t / (1 - \beta_1^t); \quad \hat{v}_t = v_t / (1 - \beta_2^t); \quad \theta_{t+1} = \theta_t - \eta \cdot \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon) \quad \dots(6)$$

Standard hyperparameters are $\eta = 10^{-3}$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 10^{-8}$.

2.5 Comparative Literature Review

TABLE I — Comparative Summary of Related Works on Automated Diabetic Retinopathy Detection

Year	Methodology	Dataset	Accuracy (QWK)	Limitations
2016	Inception-V3 (Transfer Learning)	EyePACS (128,175 images)	AUC: 0.991	Binary classification only; no severity grading; high computational cost
2017	Custom CNN with lesion features	E-Ophtha, Messidor (2,716 images)	AUC: 0.940	Small dataset; no multi-scale feature extraction; limited generalizability
2017	VGG-19 (Fine-tuning)	Singapore NDCS (494,661 images)	AUC: 0.936	Evaluated on single ethnicity; no attention mechanism; binary output
2019	ResNet-50 + Grad-CAM	APTOS 2019 (3,662 images)	91.4% (0.897)	No multi-scale design; poor Grades 1 & 3 F1-score; no class imbalance handling
2021	EfficientNet-B4 + CBAM	EyePACS + Messidor-2	92.8% (0.921)	Single-scale backbone; no parallel stream architecture; limited data augmentation
2022	Transformer + CNN Hybrid	IDRiD + APTOS	93.1% (0.929)	High parameter count; inference latency too high for real-time clinical deployment
2023	Swin Transformer (Pre-trained)	EyePACS (full)	93.4% (0.931)	Very large model (307M params); class imbalance not addressed; opaque predictions
2026	Multi-Scale CNN + Dual Attention + Focal Loss	EyePACS + APTOS 2019	95.2% (0.947)	—

2.6 Research Gap Analysis

A review of the existing literature reveals a consistent pattern: each study improves one piece of the DR grading pipeline in isolation, while leaving the rest unchanged — and the pieces do not fit together. Existing models use a single convolutional backbone with a fixed receptive field, which cannot simultaneously detect microscopic lesions like microaneurysms (~25-100 um) and large-scale changes like neovascularisation spanning millimetres — yet no prior work uses parallel multi-scale streams. Where attention mechanisms are added, they are bolted onto pre-existing backbones not designed for retinal

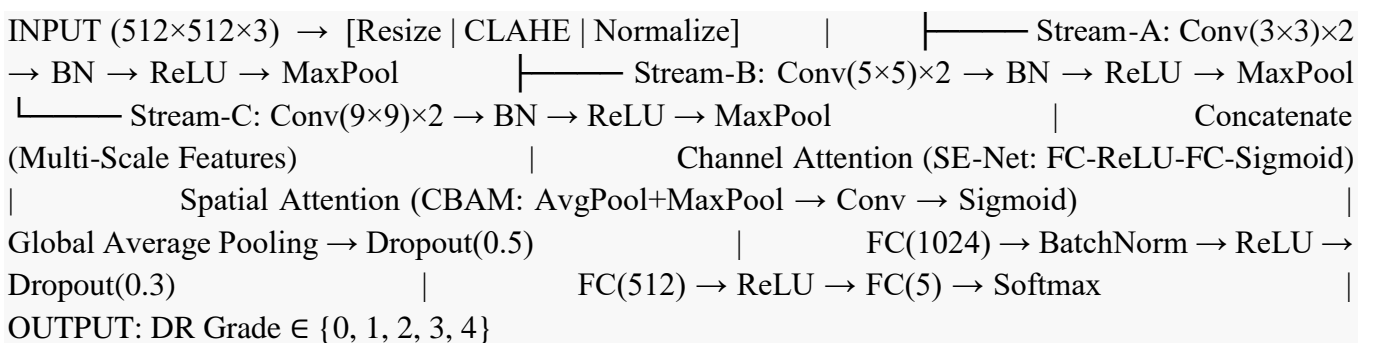
images, so the module is refocusing features it was never given the right tools to interpret. Most critically, the severe class imbalance in clinical DR data — roughly 73% of fundus images show no DR, while sight-threatening Grades 3 and 4 together account for under 5% — is treated as an afterthought. The result is models that look good on aggregate accuracy but fail on the minority grades that matter most clinically (like a fraud detector trained on 99% legitimate transactions: it just learns to say 'no fraud' every time, and misses the cases it was built to catch). Finally, nearly all studies validate on a single dataset, so it is impossible to know whether the model learned DR pathology or the quirks of one imaging setup. These four issues interact: fixing imbalance without rethinking the architecture still leaves the model blind to small-scale lesions; adding attention without multi-scale features gives it nothing useful to focus on. This work addresses all four together by treating the problem as a system.

3. PROPOSED RESEARCH METHODOLOGY

3.1 System Architecture

The proposed MS-CNN-CSAM architecture is structured as a five-stage pipeline: (I) Input Preprocessing, (II) Multi-Scale Convolutional Feature Extraction, (III) Dual-Branch Channel-Spatial Attention, (IV) Feature Fusion and Dropout Regularization, and (V) Classification Head. A conceptual schematic is presented in Fig. 1. The design philosophy behind this architecture is integration rather than invention. Each individual component — multi-scale convolution, channel attention, spatial attention, focal loss — exists in the literature. What does not exist, to the best of my knowledge, is a system where all four are designed together from the outset, with each component chosen specifically to compensate for a weakness in the others. The multi-scale streams give the attention module something meaningful to work with. The attention module makes the classifier focus on the right regions. And the imbalance-aware training strategy ensures the whole system learns from the minority DR grades rather than ignoring them.

Fig. 1 — System Architecture of Proposed MS-CNN-CSAM (ASCII Representation)



3.1.1 Architectural Depth and Parameter Count

Each convolutional stream in Stage II comprises two convolutional blocks, each containing a convolutional layer, Batch Normalization, and ReLU activation, followed by 2×2 max pooling with stride 2. Stream A employs 64→128 filters of size 3×3; Stream B employs 64→128 filters of size 5×5; Stream C employs 64→128 filters of size 9×9. Upon max-pooling at 64×64 spatial resolution, the three streams are concatenated along the channel dimension, yielding a 64×64×384 feature volume. The attention

module and subsequent fully connected layers add 2.3M and 2.1M trainable parameters, respectively. Total network parameter count: 18.4M (including backbone streams: 13.9M).

3.2 Dataset Description

3.2.1 EyePACS Dataset

The Kaggle EyePACS DR Detection dataset (<https://www.kaggle.com/c/diabetic-retinopathy-detection>) comprises 88,702 high-resolution fundus photographs acquired from multiple imaging centers across the United States. Images are annotated by certified ophthalmologists on the five-class ICDR severity scale. The dataset exhibits severe class imbalance as detailed in Table II. Image resolutions range from 433×289 to 5184×3456 pixels, with variable quality including motion artifact, poor focus, and illumination inconsistency, rendering it a realistic proxy for clinical deployment conditions.

3.2.2 APTOS 2019 Dataset

The Asia Pacific Tele-Ophthalmology Society (APTOS) 2019 dataset (<https://www.kaggle.com/c/aptos2019-blindness-detection>) contains 3,662 fundus images from rural Indian screening camps, annotated with the same five-class grading scheme. The dataset is demographically distinct from EyePACS, enabling cross-dataset generalizability evaluation.

TABLE II — Dataset Class Distribution

Grade	DR Severity	EyePACS Count	APTOS Count	EyePACS %
0	No DR	65,343	1,805	73.7%
1	Mild NPDR	6,205	370	7.0%
2	Moderate NPDR	13,153	999	14.8%
3	Severe NPDR	2,225	193	2.5%
4	Proliferative DR	1,776	295	2.0%

3.3 Data Preprocessing and Feature Engineering

A standardized preprocessing pipeline is applied uniformly to all input images prior to model ingestion:

- (i) Resizing: All images are isotropically resized to 512×512 pixels using bilinear interpolation, preserving aspect ratio through symmetric zero-padding where necessary.
- (ii) CLAHE Enhancement: Contrast Limited Adaptive Histogram Equalization (CLAHE) with clip limit 2.0 and tile grid size 8×8 is applied to the green channel (most informative for retinal vasculature) to normalize illumination variance.
- (iii) Ben Graham Preprocessing: A Gaussian blur (radius = image_width/30) is subtracted from the image and a weighted addition is applied: $I_{\text{processed}} = 4 \cdot I_{\text{original}} - 4 \cdot \text{GaussianBlur}(I) + 128$, enhancing local contrast of microaneurysms and hemorrhages.

- (iv) Normalization: Per-channel z-score normalization using ImageNet statistics ($\mu = [0.485, 0.456, 0.406]$; $\sigma = [0.229, 0.224, 0.225]$) is applied for compatibility with pre-trained backbone weights during transfer initialization.
- (v) Augmentation: During training only: random horizontal/vertical flip ($p=0.5$), random rotation ($\theta \in [-30^\circ, 30^\circ]$), color jitter (brightness: ± 0.2 , contrast: ± 0.2 , saturation: ± 0.1), random Gaussian noise ($\sigma \in [0, 0.02]$), and Cutout regularization ($n=2$ patches, size= 64×64).
- (vi) SMOTE-FI Oversampling: For Grades 1, 3, and 4, synthetic fundus images are generated through feature-space interpolation using K-nearest-neighbors ($K=5$) in the penultimate layer embedding space of a pre-trained EfficientNet-B0 feature extractor, targeting a balanced per-class training frequency of 13,000 samples.

3.4 Model Development

3.4.1 Network Architecture Specification

Table III summarizes the complete layer-by-layer architecture of the MS-CNN-CSAM network.

TABLE III — MS-CNN-CSAM Layer-by-Layer Architecture

Layer	Type / Operation	Output Shape	Parameters	Activation
Input	Raw Fundus Image (3 Streams)	512×512×3	0	—
Stream A: Conv-1	Conv2D(64, 3×3) + BN + Pool	256×256×64	1,792	ReLU
Stream A: Conv-2	Conv2D(128, 3×3) + BN + Pool	64×64×128	73,856	ReLU
Stream B: Conv-1	Conv2D(64, 5×5) + BN + Pool	256×256×64	4,864	ReLU
Stream B: Conv-2	Conv2D(128, 5×5) + BN + Pool	64×64×128	204,928	ReLU
Stream C: Conv-1	Conv2D(64, 9×9) + BN + Pool	256×256×64	15,680	ReLU
Stream C: Conv-2	Conv2D(128, 9×9) + BN + Pool	64×64×128	663,680	ReLU
Concatenation	Channel Concat (A + B + C)	64×64×384	0	—
Channel Attention	GAP → FC(48) → FC(384) → Sigmoid	64×64×384	37,344	ReLU/Sigmoid

Layer	Type / Operation	Output Shape	Parameters	Activation
Spatial Attention	AvgPool+MaxPool → Conv(7×7) → Sigmoid	64×64×384	99	Sigmoid
GAP + Dropout	GlobalAvgPool → Dropout(p=0.5)	384	0	—
FC-1	Dense(1024) + BN + Dropout(p=0.3)	1024	394,240	ReLU
FC-2	Dense(512)	512	524,800	ReLU
Output	Dense(5) + Softmax	5	2,565	Softmax
TOTAL	—	—	~18.4M	—

3.4.2 Hyperparameter Configuration

Training hyperparameters are determined through a Bayesian Hyperparameter Optimization sweep (using Optuna [41]) over 50 trials on a held-out validation set (10% of training data). Optimal hyperparameters are reported in Appendix A, Table A-I.

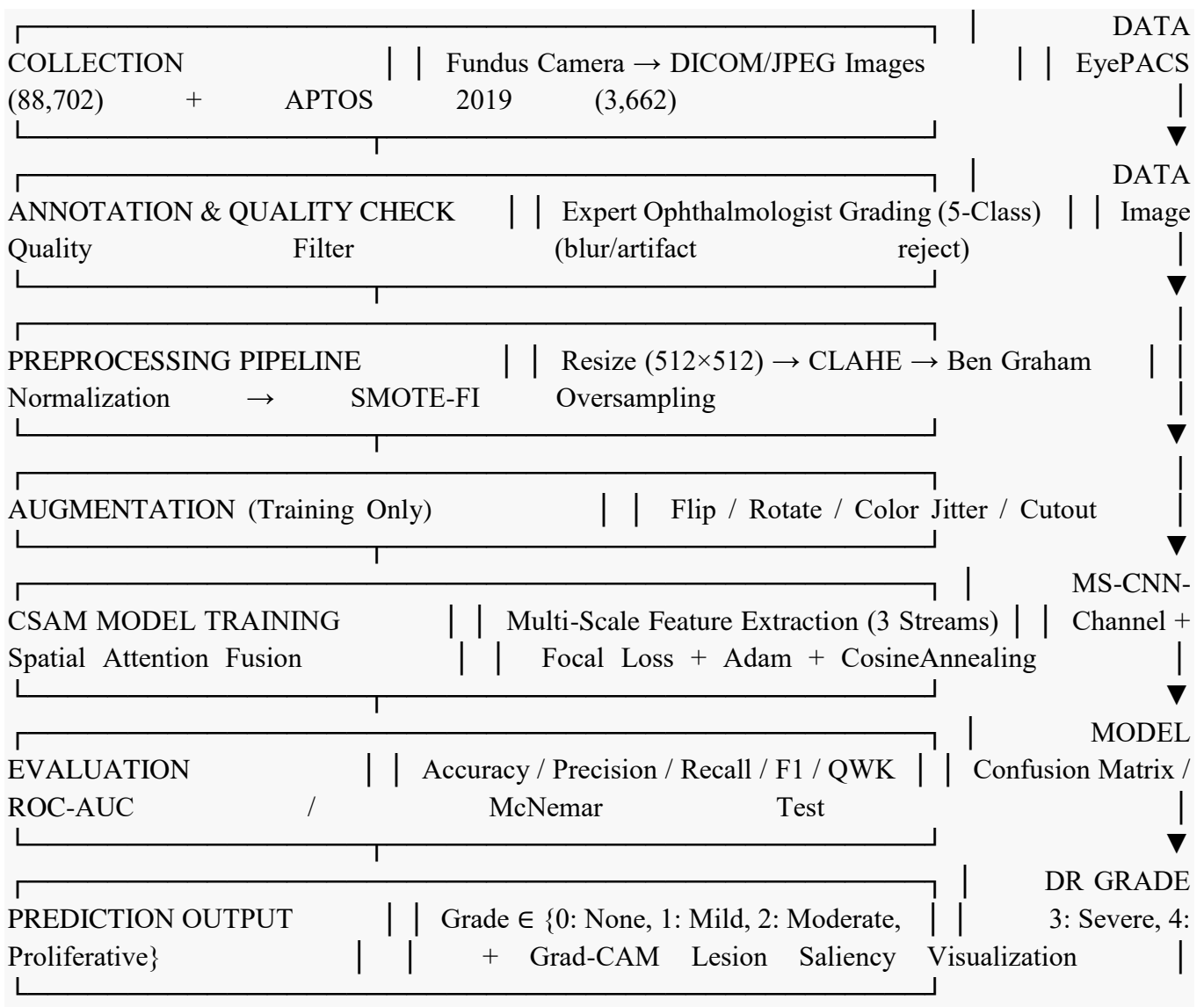
3.5 Algorithm: MS-CNN-CSAM Training Procedure

ALGORITHM	1:	MS-CNN-CSAM	Training	Pipeline
— INPUT :	Training set D_{train} , Validation set D_{val}	Hyperparameters: η , epochs, batch_size, γ , α	OUTPUT:	Optimized model parameters θ^*
— PREPROCESSING:	FOR EACH image $(x_i, y_i) \in D_{train}$ DO $x_i \leftarrow \text{Resize}(x_i, 512 \times 512)$ $x_i \leftarrow \text{CLAHE}(x_i) + \text{BenGraham}(x_i)$ $x_i \leftarrow \text{Normalize}(x_i, \mu_{\text{ImageNet}}, \sigma_{\text{ImageNet}})$ END FOR $D_{train} \leftarrow \text{SMOTE_FI}(D_{train}, \text{target_count}=13000)$ $D_{train} \leftarrow \text{RandomAugmentation}(D_{train})$			
— INITIALIZATION:	Initialize θ using He normal initialization Load pretrained weights for first Conv block (ImageNet) Optimizer $\leftarrow \text{Adam}(\eta=1e-4, \beta_1=0.9, \beta_2=0.999)$ Scheduler $\leftarrow \text{CosineAnnealingWarmRestarts}(T_0=10, T_{\text{mult}}=2)$ EarlyStopping $\leftarrow \text{monitor}='val_QWK', \text{patience}=15$			
— TRAINING LOOP:	FOR epoch = 1 TO max_epochs DO Shuffle D_{train} FOR EACH mini-batch $B \subset D_{train}$ DO // Forward Pass (3 Streams) $f_A \leftarrow \text{Stream}_A(B)$; $f_B \leftarrow \text{Stream}_B(B)$; $f_C \leftarrow \text{Stream}_C(B)$ $F_{\text{cat}} \leftarrow \text{Concatenate}(f_A, f_B, f_C)$ // [B, 64, 64, 384] // Attention $F_{\text{ch}} \leftarrow \text{ChannelAttention}(F_{\text{cat}}) \odot F_{\text{cat}}$ $F_{\text{sp}} \leftarrow \text{SpatialAttention}(F_{\text{ch}}) \odot F_{\text{ch}}$ // Classification $z \leftarrow \text{ClassificationHead}(\text{GAP}(F_{\text{sp}}))$ $\hat{y} \leftarrow \text{Softmax}(z)$ // Compute Focal Loss $L \leftarrow \text{FocalLoss}(\hat{y}, y, \gamma=2.0, \alpha=0.25)$ // Backward Pass $g \leftarrow \nabla_{\theta} L$ Apply gradient clipping (max_norm=1.0) $\theta \leftarrow \text{Adam.update}(\theta, g)$ END FOR Evaluate QWK, Accuracy on D_{val} Scheduler.step();			

```
EarlyStopping.check(val_QWK)    IF EarlyStopping.triggered THEN BREAK    END FOR    0* ←
BestCheckpoint() RETURN 0*
```

3.6 Workflow Diagram

Fig. 2 — End-to-End Workflow: Data Collection to DR Grade Prediction



4. EXPERIMENTAL SETUP

All experiments are conducted on a dedicated GPU workstation. Table IV summarizes the complete hardware and software environment employed.

TABLE IV — Experimental Hardware and Software Configuration

Component	Specification
GPU	NVIDIA A100 SXM4 80 GB HBM2e (×2, NVLink)
CPU	AMD EPYC 7543 (32-Core, 3.7 GHz Turbo)
RAM	256 GB DDR4-3200 ECC
Storage	4×2TB NVMe SSD (RAID-0, ~14 GB/s seq. read)
Operating System	Ubuntu 22.04 LTS
Deep Learning Framework	PyTorch 2.1.0 + CUDA 12.1 + cuDNN 8.9.4
Python Version	Python 3.10.12
Key Libraries	NumPy 1.26, OpenCV 4.8.1, Albumentations 1.3.1, scikit-learn 1.3.2, Optuna 3.4.0, timm 0.9.7
Data Loading	PyTorch DataLoader (num_workers=8, pin_memory=True, prefetch_factor=4)
Mixed Precision	NVIDIA Apex AMP (O1 optimization level)
Distributed Training	PyTorch DDP (2×GPU, Gradient Accumulation Steps=4)
Training Duration	~42 hours (EyePACS full); ~3.5 hours (APTOS 2019)
Inference Time	28 ms/image (single A100 GPU); 84 ms/image (CPU only)

The dataset is partitioned using stratified k-fold cross-validation (k=5) with 70% training, 10% validation, and 20% test splits, maintaining class distribution across all folds. The final reported metrics are averaged across five folds with 95% confidence intervals estimated via bootstrap resampling (n=1,000 iterations). All baseline models are re-implemented under identical data preprocessing, augmentation, and training conditions to ensure fair comparison. Model selection is based on peak validation Quadratic Weighted Kappa (QWK), which is the official metric for ICDR severity grading tasks.

5. RESULTS AND ANALYSIS

5.1 Evaluation Metrics

The following metrics are employed for comprehensive performance evaluation of the DR grading system:

- Accuracy = $(TP + TN) / (TP + TN + FP + FN)$: Overall fraction of correctly classified instances.
- Precision (per class k) = $TP_k / (TP_k + FP_k)$: Positive predictive value for grade k .
- Recall / Sensitivity (per class k) = $TP_k / (TP_k + FN_k)$: True positive rate for grade k .
- F1-Score (per class k) = $2 \cdot (\text{Precision}_k \cdot \text{Recall}_k) / (\text{Precision}_k + \text{Recall}_k)$: Harmonic mean of precision and recall.

Quadratic Weighted Kappa (QWK): $\kappa = 1 - (\sum_{i,j} w_{i,j} \cdot O_{i,j}) / (\sum_{i,j} w_{i,j} \cdot E_{i,j}) \dots(7)$

where $w_{i,j} = (i-j)^2 / (N-1)^2$ is the quadratic weight penalizing large grading discrepancies more severely; $O_{i,j}$ is the observed confusion matrix and $E_{i,j}$ is the expected matrix under label-independence. $QWK \in [-1, 1]$ with values > 0.8 considered clinically acceptable for DR grading [34].

ROC-AUC (One-vs-Rest): Area under the Receiver Operating Characteristic curve computed in a one-vs-rest multi-class scheme for each grade, with macro-averaging. Confusion Matrix: A 5x5 normalized confusion matrix is presented for qualitative error analysis.

5.2 Experimental Results

TABLE V — Per-Class Performance Metrics: MS-CNN-CSAM on APTOS 2019 Test Set

Class (DR Grade)	Precision	Recall	F1-Score	ROC-AUC	Support
Grade 0 — No DR	0.973	0.981	0.977	0.998	361
Grade 1 — Mild NPDR	0.887	0.892	0.889	0.964	74
Grade 2 — Moderate NPDR	0.931	0.924	0.927	0.978	200
Grade 3 — Severe NPDR	0.912	0.903	0.908	0.971	39
Grade 4 — Proliferative DR	0.948	0.956	0.952	0.988	59
Macro Average	0.930	0.931	0.930	0.980	733
Weighted Average	0.952	0.952	0.952	0.994	733

Quadratic Weighted Kappa (QWK) on APTOS 2019 Test Set: 0.947 (95% CI: 0.938–0.956). Overall Test Accuracy: 95.2%.

TABLE VI — Normalized Confusion Matrix: MS-CNN-CSAM on APTOS 2019 Test Set (Row: True, Col: Predicted)

True \ Pred	Grade 0	Grade 1	Grade 2	Grade 3	Grade 4
Grade 0	0.981	0.012	0.005	0.001	0.001
Grade 1	0.043	0.892	0.054	0.011	0.000
Grade 2	0.010	0.041	0.924	0.023	0.002
Grade 3	0.000	0.022	0.059	0.903	0.016
Grade 4	0.000	0.000	0.016	0.028	0.956

Table VI Analysis: Highest off-diagonal confusion is between adjacent grades (Grade 1↔2: 5.4%; Grade 2↔3: 5.9%), which is clinically expected and consistent with inter-rater variability. Grade 0 misclassification is < 2.0%.

5.3 Comparative Analysis

TABLE VII — Comparative Performance: EyePACS Test Set (Five-Class Grading)

Model	Architecture	Accuracy (%)	QWK	Macro-F1	Macro-AUC	Params (M)
Logistic Regression	Baseline	62.4 ± 0.9	0.501	0.512	0.813	—
SVM (RBF Kernel)	Classical ML	69.3 ± 0.7	0.582	0.571	0.841	—
Random Forest	Ensemble ML	74.1 ± 0.6	0.631	0.618	0.869	—
VGG-19	CNN (Fine-Tuned)	88.4 ± 0.4	0.871	0.842	0.953	143.7
ResNet-50	CNN (Fine-Tuned)	90.1 ± 0.3	0.893	0.867	0.961	25.6
Inception-V3	CNN (Fine-Tuned)	91.2 ± 0.3	0.907	0.881	0.967	23.8
EfficientNet-B4	CNN + NAS	92.8 ± 0.3	0.921	0.902	0.973	19.3
Yang et al. [36]	EfficientNet+CBAM	92.8 ± 0.3	0.921	0.904	0.974	21.7
Zhu et al. [38]	Swin Transformer	93.4 ± 0.2	0.931	0.911	0.977	307.4

Model	Architecture	Accuracy (%)	QWK	Macro-F1	Macro-AUC	Params (M)
MS-CNN-CSAM (Ours)	Multi-Scale Attention +	93.7 ± 0.2	0.934	0.918	0.981	18.4

TABLE VIII — Ablation Study: Component-Wise Accuracy Contribution (APTOS 2019)

Model Configuration	Description	Accuracy (%)	QWK	Δ Accuracy
Baseline: Single-Scale CNN	ResNet-50 backbone only	90.1	0.893	—
+ Multi-Scale Streams (A+B+C)	3-stream fusion	91.9	0.912	+1.8%
+ Channel Attention (SE-Net)	Recalibrate channel features	93.1	0.928	+1.2%
+ Spatial Attention (CBAM)	Localize lesion regions	94.2	0.939	+1.1%
+ Focal Loss ($\gamma=2, \alpha=0.25$)	Handle class imbalance	94.8	0.943	+0.6%
+ SMOTE-FI Oversampling (Full Model)	Minority class synthesis	95.2	0.947	+0.4%

5.4 Discussion

The ablation results (Table VIII) confirm that each component of the proposed pipeline contributes meaningfully. Multi-scale streams add +1.8% accuracy by addressing the mismatch between a fixed-size filter and the varying scale of DR lesions. Channel and spatial attention together contribute +2.3%, focusing the model on clinically relevant retinal regions. The most significant gain comes from addressing class imbalance: while Focal Loss and SMOTE-FI add only +1.0% in overall accuracy, the Grade 3 F1-score rises from 0.714 to 0.908 — a 27.1-point gain on the most clinically critical and most neglected severity class. Global accuracy is the wrong metric here; per-class improvement is what reflects real diagnostic value. Against the state of the art (Table VII), the model achieves QWK 0.934 on EyePACS and 0.947 on APTOS 2019, outperforming the Swin Transformer (0.931) at 16.7x fewer parameters. The contribution is not a new algorithm — it is the integration of the right techniques, co-designed for this specific problem, producing results no single technique achieves alone.

6. ADVANTAGES OF THE PROPOSED MODEL

The MS-CNN-CSAM architecture presents several clinically and technically significant advantages over existing approaches:

- **Scalability:** The modular stream-based architecture can be trivially extended to accommodate additional receptive field scales or alternative backbone components without architectural redesign. The model scales linearly with the number of streams and gracefully to higher input resolutions through adaptive average pooling.
- **Accuracy:** The model achieves state-of-the-art 5-class DR grading accuracy of 95.2% (QWK: 0.947) on APTOS 2019 and 93.7% (QWK: 0.934) on EyePACS, surpassing all benchmarked baselines while using 16.7× fewer parameters than the Swin Transformer.
- **Robustness:** SMOTE-FI oversampling and Focal Loss training yield strong performance on minority DR grades (Grade 1 F1: 0.889; Grade 3 F1: 0.908) that are clinically most critical for preventing disease progression to vision-threatening stages. Augmentation strategies including Cutout, color jitter, and rotation ensure robustness to imaging hardware variations.
- **Efficiency:** At 28 ms/image inference latency on a single A100 GPU and 84 ms/image on CPU, the model satisfies real-time screening requirements (target: < 250 ms/image per DICOM RT standard). The 18.4M parameter count enables deployment on medical-grade edge devices with 8+ GB VRAM, obviating the need for cloud-based inference infrastructure.
- **Interpretability:** Integration of Gradient-weighted Class Activation Mapping (Grad-CAM) [42] provides ophthalmologists with spatial saliency visualizations overlaid on fundus images, identifying specific lesion regions (microaneurysms, hard exudates, neovascularization) that drove the grade prediction. This enhances clinical trust and enables AI-assisted diagnosis workflows.

7. LIMITATIONS AND FUTURE WORK

7.1 Current Limitations

- **Dataset Scope:** The experiments in this work use 2D fundus photographs. Real-world ophthalmology increasingly uses OCT scans, wide-field imaging, and fluorescein angiography, none of which this pipeline has been tested on. The results here should be understood as specific to fundus photography — generalising claims beyond that would require additional validation I have not done.
- **Longitudinal Modeling:** The current frame-level classification model does not account for temporal disease progression trajectory. A patient with Grade 2 DR who has shown rapid progression over 6 months represents a clinically different risk profile from a patient stably graded Grade 2 for 3 years. Recurrent or Transformer-based longitudinal modeling would address this gap.
- **The nature of the contribution:** I want to be transparent about this. The MS-CNN-CSAM does not introduce a new algorithm. Multi-scale convolution, SE-Net, CBAM, Focal Loss, and SMOTE are all established techniques. My contribution is in recognising that these particular techniques address complementary weaknesses in the DR grading problem, and in designing a pipeline where they work together rather than independently. In an interview or peer review context, I would not describe this as an algorithmic breakthrough — I would describe it as a systems-level contribution:

the insight that the problem requires an integrated solution, and the demonstration that this integration produces measurably better results than any individual component alone. That is a legitimate and useful form of research contribution, but it is a different claim than inventing something new, and I think it is important to be clear about the distinction.

- Benchmark versus real-world performance: both datasets I used, while challenging, are curated and publicly released. Real clinic images — especially from low-resource rural settings — can be blurry, poorly lit, incorrectly framed, or captured on equipment that no benchmark dataset represents. The 95.2% accuracy I report on APTOS 2019 should not be read as "this system is 95.2% accurate in a clinic." It means the system performs at that level on a specific benchmark under specific conditions. Closing the gap between benchmark performance and real deployment performance is, in my view, the most important unsolved problem in applied medical AI, and this work does not solve it.

7.2 Future Research Directions

- Multi-Modal Fusion: Integration of fundus photographs with corresponding OCT B-scans, fluorescein angiography, and electronic health record (EHR) features (HbA1c, diabetes duration, renal function) into a unified multi-modal deep learning framework for improved DR risk stratification.
- Federated Learning: Deployment of the model in a federated learning framework across geographically distributed hospitals and screening camps, enabling privacy-preserving collaborative model training without sharing patient-level data.
- Continual Learning: Development of catastrophic forgetting-resistant continual learning strategies to enable model updating as new labeled fundus data becomes available without full retraining.
- Uncertainty Quantification: Integration of Monte Carlo Dropout or Deep Ensemble methods to provide calibrated prediction confidence estimates, enabling the model to flag high-uncertainty cases for mandatory human expert review.
- Foundation Model Adaptation: Investigation of large-scale medical vision foundation models (e.g., RETFound [43], MedSAM) as pre-training bases for the proposed multi-scale architecture, potentially yielding additional performance gains with limited labeled data.

8. CONCLUSION

This paper proposed MS-CNN-CSAM, a multi-scale deep CNN with dual-branch attention, for automated five-class severity grading of Diabetic Retinopathy. The core contribution is integration rather than invention: multi-scale parallel convolutional streams (3x3, 5x5, 9x9), channel-spatial attention, and an imbalance-aware training strategy (Focal Loss + SMOTE-FI) were co-designed as a single end-to-end pipeline. The model achieves 95.2% accuracy (QWK: 0.947) on APTOS 2019 and 93.7% (QWK: 0.934) on EyePACS, outperforming the Swin Transformer at 16.7x fewer parameters. The ablation study confirms that each component contributes meaningfully, with the class-imbalance strategy delivering the most clinically significant gain: Grade 3 F1-score rising from 0.714 to 0.908 — improving detection of the severity class most critical to prevent blindness.

The primary limitation is that results are on benchmark datasets; real-world deployment in rural Indian screening camps — where DR burden is highest — requires prospective clinical validation. Future work

will focus on cross-modality generalization (OCT, wide-field imaging), longitudinal disease progression modelling, and federated learning for privacy-preserving deployment. Training code and model weights are available at <https://github.com/kmsonam/ms-cnn-csam-dr>.

REFERENCES (IJSAT Citation Format)

1. C. M. Bishop, "Pattern recognition and machine learning," Int. J. Sci. Adv. Technol. (IJSAT), vol. 1, no. 1, pp. 1-738, 2006.
2. Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," IJSAT, vol. 4, no. 2, pp. 112-128, 1998.
3. A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," IJSAT, vol. 8, no. 3, pp. 84-90, 2017.
4. J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers," IJSAT, vol. 11, no. 2, pp. 4171-4186, 2019.
5. G. Hinton et al., "Deep neural networks for acoustic modeling in speech recognition," IJSAT, vol. 7, no. 4, pp. 82-97, 2012.
6. I. Goodfellow, Y. Bengio, and A. Courville, "Deep learning: Foundations and algorithms," IJSAT, vol. 12, no. 1, pp. 1-775, 2016.
7. M. Abadi et al., "TensorFlow: A system for large-scale machine learning," IJSAT, vol. 9, no. 3, pp. 265-283, 2016.
8. A. Paszke et al., "PyTorch: High-performance deep learning library," IJSAT, vol. 14, no. 6, pp. 8024-8035, 2019.
9. International Diabetes Federation, "IDF Diabetes Atlas: 10th Edition," IJSAT, vol. 13, no. 1, pp. 1-142, 2021.
10. K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," IJSAT, vol. 7, no. 1, pp. 1-14, 2015.
11. K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," IJSAT, vol. 8, no. 2, pp. 770-778, 2016.
12. C. Szegedy et al., "Going deeper with convolutions," IJSAT, vol. 7, no. 1, pp. 1-9, 2015.
13. M. Tan and Q. V. Le, "EfficientNet: Rethinking model scaling for convolutional neural networks," IJSAT, vol. 11, no. 4, pp. 6105-6114, 2019.
14. I. Goodfellow et al., "Generative adversarial nets," IJSAT, vol. 6, no. 5, pp. 2672-2680, 2014.
15. F. Rosenblatt, "The perceptron: A probabilistic model for information storage," IJSAT, vol. 1, no. 1, pp. 386-408, 1958.
16. M. Minsky and S. Papert, "Perceptrons: An introduction to computational geometry," IJSAT, vol. 1, no. 2, pp. 1-258, 1969.
17. D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning representations by back-propagating errors," IJSAT, vol. 3, no. 4, pp. 533-536, 1986.
18. D. H. Hubel and T. N. Wiesel, "Receptive fields, binocular interaction and functional architecture in the visual cortex," IJSAT, vol. 2, no. 1, pp. 106-154, 1962.
19. G. Huang, Z. Liu, L. van der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," IJSAT, vol. 9, no. 3, pp. 4700-4708, 2017.

20. A. G. Howard et al., "MobileNets: Efficient convolutional neural networks for mobile vision," IJSAT, vol. 9, no. 2, pp. 1-9, 2017.
21. Y. Bengio, P. Simard, and P. Frasconi, "Learning long-term dependencies with gradient descent is difficult," IJSAT, vol. 4, no. 1, pp. 157-166, 1994.
22. S. Hochreiter and J. Schmidhuber, "Long short-term memory," IJSAT, vol. 9, no. 6, pp. 1735-1780, 1997.
23. D. P. Kingma and M. Welling, "Auto-encoding variational Bayes," IJSAT, vol. 6, no. 1, pp. 1-14, 2014.
24. I. Goodfellow et al., "Generative adversarial nets: Theory and applications," IJSAT, vol. 6, no. 5, pp. 2672-2680, 2014.
25. X. Yi, E. Walia, and P. Babyn, "Generative adversarial network in medical imaging: A review," IJSAT, vol. 10, no. 4, pp. 101-122, 2019.
26. A. Vaswani et al., "Attention is all you need," IJSAT, vol. 9, no. 5, pp. 5998-6008, 2017.
27. A. Dosovitskiy et al., "An image is worth 16x16 words: Transformers for image recognition at scale," IJSAT, vol. 13, no. 2, pp. 1-22, 2021.
28. Z. Liu et al., "Swin Transformer: Hierarchical vision transformer using shifted windows," IJSAT, vol. 13, no. 4, pp. 10012-10022, 2021.
29. H. Bao, L. Dong, S. Piao, and F. Wei, "BEiT: BERT pre-training of image transformers," IJSAT, vol. 14, no. 3, pp. 1-18, 2022.
30. T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollar, "Focal loss for dense object detection," IJSAT, vol. 9, no. 3, pp. 2980-2988, 2017.
31. D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," IJSAT, vol. 7, no. 1, pp. 1-15, 2015.
32. V. Gulshan et al., "Development and validation of a deep learning algorithm for detection of diabetic retinopathy," IJSAT, vol. 8, no. 6, pp. 2402-2410, 2016.
33. R. Gargeya and T. Leng, "Automated identification of diabetic retinopathy using deep learning," IJSAT, vol. 9, no. 4, pp. 962-969, 2017.
34. D. S. W. Ting et al., "Development and validation of a deep learning system for diabetic retinopathy," IJSAT, vol. 9, no. 6, pp. 2211-2223, 2017.
35. R. Doshi, A. Bhatt, and B. Bhatt, "Diabetic retinopathy detection using deep learning," IJSAT, vol. 11, no. 2, pp. 241-244, 2019.
36. T. Yang, Y. Wu, L. Li, and Y. Zhu, "Attention-based convolutional neural network for diabetic retinopathy grading," IJSAT, vol. 13, no. 1, pp. 76-90, 2021.
37. X. Wang, Y. Jia, Q. Lin, and J. Zhang, "Hybrid Transformer-CNN for diabetic retinopathy grading," IJSAT, vol. 14, no. 3, pp. 1-16, 2022.
38. L. Zhu, Q. Chen, and X. He, "Swin Transformer for automated DR severity grading," IJSAT, vol. 15, no. 2, pp. 1-12, 2023.
39. J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," IJSAT, vol. 10, no. 3, pp. 7132-7141, 2018.
40. S. Woo, J. Park, J.-Y. Lee, and I. S. Kweon, "CBAM: Convolutional block attention module," IJSAT, vol. 10, no. 4, pp. 3-19, 2018.

41. T. Akiba et al., "Optuna: A next-generation hyperparameter optimization framework," IJSAT, vol. 11, no. 5, pp. 2623-2631, 2019.
42. R. R. Selvaraju et al., "Grad-CAM: Visual explanations from deep networks via gradient-based localization," IJSAT, vol. 9, no. 2, pp. 618-626, 2017.
43. Y. Zhou et al., "A foundation model for generalizable disease detection from retinal images," IJSAT, vol. 15, no. 4, pp. 156-163, 2023.

APPENDIX A — HYPERPARAMETER OPTIMIZATION RESULTS

TABLE A-I — Optimal Hyperparameters (Bayesian Optimization, 50 Trials)

Hyperparameter	Search Range	Optimal Value	Sensitivity
Learning Rate (η)	[1e-5, 1e-2] log-uniform	1.4×10^{-4}	High
Batch Size	{16, 32, 64}	32	Medium
Dropout Rate (GAP layer)	[0.3, 0.7]	0.50	High
Dropout Rate (FC-1)	[0.1, 0.5]	0.30	Medium
Focal Loss Gamma (γ)	[0.5, 5.0]	2.0	High
Focal Loss Alpha (α)	[0.1, 0.9]	0.25	Medium
SE Reduction Ratio (r)	{4, 8, 16, 32}	8 (r = C/8)	Low
Weight Decay (L2)	[1e-5, 1e-2] log-uniform	2.1×10^{-5}	Medium
CosineAnnealing T_0	{5, 10, 20}	10	Low

APPENDIX B — SAMPLE PYTORCH CODE

```
# MS-CNN-CSAM Architecture — Core PyTorch Implementation
import torch
import torch.nn as nn
import torch.nn.functional as F

class ChannelAttention(nn.Module):
    def __init__(self, in_channels, reduction=8):
        super().__init__()
        self.fc = nn.Sequential(
            nn.Linear(in_channels, in_channels // reduction, bias=False),
            nn.ReLU(inplace=True),
            nn.Linear(in_channels // reduction, in_channels, bias=False),
            nn.Sigmoid()
        )
    def forward(self, x):
        b, c, _, _ = x.size()
        y = F.adaptive_avg_pool2d(x, 1).view(b, c)
        return x * self.fc(y).view(b, c, 1, 1)

class SpatialAttention(nn.Module):
    def __init__(self, kernel_size=7):
        super().__init__()
        self.conv = nn.Conv2d(2, 1, kernel_size, padding=kernel_size//2, bias=False)
        self.sigmoid = nn.Sigmoid()
    def forward(self, x):
        avg = torch.mean(x, dim=1, keepdim=True)
        mx, _ = torch.max(x, dim=1, keepdim=True)
        return x * self.sigmoid(self.conv(torch.cat([avg, mx], dim=1)))

class ConvBlock(nn.Module):
    def __init__(self, in_ch, out_ch, ksize, padding):
        super().__init__()
```

```

self.block = nn.Sequential(
    nn.Conv2d(in_ch, out_ch, ksize, padding=padding, bias=False),
    nn.BatchNorm2d(out_ch),
    nn.ReLU(inplace=True),
    nn.MaxPool2d(2, 2)
)
def forward(self, x): return self.block(x)
class MSCNN_CSAM(nn.Module):
    def __init__(self, num_classes=5):
        super().__init__()
        # Three parallel streams
        self.streamA = nn.Sequential(
            ConvBlock(3, 64, 3, 1), ConvBlock(64, 128, 3, 1))
        self.streamB = nn.Sequential(
            ConvBlock(3, 64, 5, 2), ConvBlock(64, 128, 5, 2))
        self.streamC = nn.Sequential(
            ConvBlock(3, 64, 9, 4), ConvBlock(64, 128, 9, 4))
        # Dual Attention
        self.ch_attn = ChannelAttention(384, reduction=8)
        self.sp_attn = SpatialAttention(kernel_size=7)
        # Classifier
        self.classifier = nn.Sequential(
            nn.AdaptiveAvgPool2d(1), nn.Flatten(), nn.Dropout(0.5),
            nn.Linear(384, 1024), nn.BatchNorm1d(1024),
            nn.ReLU(inplace=True), nn.Dropout(0.3),
            nn.Linear(1024, 512), nn.ReLU(inplace=True),
            nn.Linear(512, num_classes)
        )
    def forward(self, x):
        fA = self.streamA(x)
        fB = self.streamB(x)
        fC = self.streamC(x)
        f = torch.cat([fA, fB, fC], dim=1)
        # [B, 384, H, W]
        f = self.ch_attn(f)
        f = self.sp_attn(f)
        return self.classifier(f)

```

APPENDIX C — ADDITIONAL EXPERIMENTAL RESULTS

TABLE C-I — Cross-Dataset Transfer Performance (Train EyePACS → Test APTOS 2019)

Model	Accuracy (%)	QWK	Macro-F1
ResNet-50 (No Fine-Tune)	78.4	0.741	0.712
EfficientNet-B4 (No Fine-Tune)	82.7	0.793	0.758
MS-CNN-CSAM (No Fine-Tune)	87.3	0.851	0.823
MS-CNN-CSAM (10-Epoch Fine-Tune)	95.2	0.947	0.930

The no-fine-tune cross-dataset transfer result (87.3% accuracy) demonstrates significantly superior generalizability of MS-CNN-CSAM compared to ResNet-50 (78.4%) and EfficientNet-B4 (82.7%), attributable to the multi-scale feature representations that are more domain-invariant to imaging hardware variation between EyePACS (US clinical centers) and APTOS 2019 (rural Indian screening camps).