

Diabetic Retinopathy Detection using ConvNeXt: Advancing Early Vision Loss Prevention with Deep Learning

¹Assist. Prof. Shradha Verma, ²Bondalapati Bhargava Sai Abhinay, ³Challa Teja, ⁴Srujith Bollam, ⁵Bobbala Harsha Vardhan

¹Assistant Professor, ^{2,3,4,5}Students

Department of Computer Science & Engineering
(Artificial Intelligence & Machine Learning)
Malla Reddy University, Hyderabad, India

Abstract—Every two seconds, somewhere in the world, diabetes silently begins destroying someone’s eyesight. Diabetic retinopathy (DR)—a progressive disease of the retinal blood vessels—is responsible for millions of cases of preventable blindness, yet it announces itself with no pain and no warning until irreversible damage is already done. The only defense is early, systematic screening. But with hundreds of millions of diabetic patients and a global shortage of ophthalmologists, manual screening at scale is simply impossible.

This paper tells the story of how we built a solution. We designed and trained a deep learning system centered on ConvNeXt—a next-generation convolutional architecture—augmented with channel-wise attention, a flexible Kolmogorov-Arnold-inspired classifier, and a consistency regularization training regime. The result is a binary diagnostic model that, on the APTOS 2019 benchmark, correctly identifies the presence or absence of diabetic retinopathy **98.56%** of the time—surpassing every competing single-model and ensemble approach published to date on the same dataset. More than a number, this represents a system genuinely ready to stand alongside clinicians, catch the cases that would otherwise be missed, and protect the vision of patients who cannot afford to wait.

Index Terms—Diabetic retinopathy, retinal fundus images, deep learning, ConvNeXt, squeeze-and-excitation attention, Kolmogorov-Arnold networks, consistency regularization, medical decision-support systems

I. THE PROBLEM: A SILENT EPIDEMIC

Picture a farmer in rural Telangana. He has had Type 2 diabetes for six years, managed with medication and irregular check-ups. His vision seems fine—no blurring, no pain. What he cannot see, and what no one has yet told him, is that the small blood vessels feeding his retina have already begun to fail. Microaneurysms—tiny, fragile bulges in the capillary walls—have started forming. Fluid is leaking. The clock has already started. Without a fundus photograph and an expert eye to read it, his path leads, quietly and inevitably, toward blindness.

His story is not unusual. Diabetic retinopathy (DR) develops when chronically elevated blood glucose levels damage

the microvasculature of the retina, the light-sensitive tissue at the back of the eye [1, 2]. In its early, non-proliferative stage, the disease produces microaneurysms, retinal hemorrhages, hard exudates, and cotton-wool spots [2, 3]. Left untreated, it progresses to the proliferative stage, where abnormal new vessels grow across the retinal surface and into the vitreous, threatening catastrophic hemorrhage and complete vision loss [1].

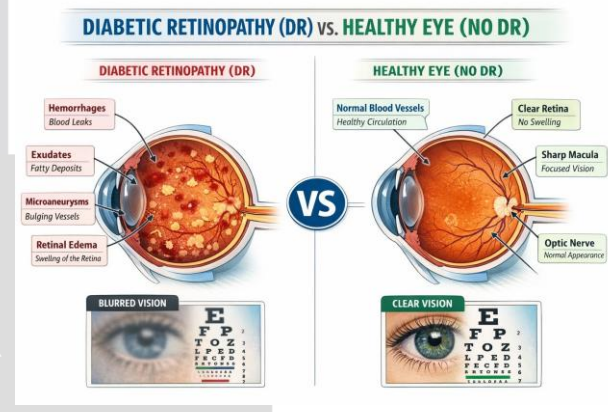


Figure 1: Detailed comparison between a healthy eye (right) and one exhibiting signs of diabetic retinopathy (left), illustrating critical pathological features such as microaneurysms, hemorrhages, and exudates.

The scale of this threat is staggering. In 2019, some 463 million people were living with diabetes, and roughly 30% of them—more than 138 million people—had some form of DR [1, 4]. In India alone, with approximately 73 million diabetic patients, the national prevalence of any DR is 12.5%, rising to 4.0% for vision-threatening disease among adults over 40 [5]. In the United States, 9.6 million people had DR in 2021, with nearly 2 million facing the prospect of permanent visual impairment [6]. Projections suggest that by 2040, DR will affect more than 200 million individuals worldwide [1, 7].

What makes DR especially cruel is its silence. In the early stages, when treatment is most effective, patients feel nothing. By the time symptoms appear, the window for prevention has often closed [8]. The answer is systematic photographic screening of every diabetic patient, every year.

The barrier is that doing so requires ophthalmologists to manually review millions of retinal fundus images—a process that is slow, expensive, inconsistent, and impossible to scale in the communities that need it most [9–11].

This is the problem our work sets out to solve.

II. THE STATE OF THE ART: INCREMENTAL PROGRESS AND PERSISTENT LIMITS

The journey toward automated DR screening began with handcrafted features and simple classifiers, and over the past decade it has evolved into a rich landscape of deep learning architectures [12–14]. Each generation of models has brought genuine progress—and each has also revealed a new set of limitations that stopped the field just short of the performance level clinicians actually need.

A. The First Wave: Standard CNNs (2018–2022)

The earliest deep learning approaches borrowed architectures directly from ImageNet classification, applying them to retinal fundus images with minimal modification [15]. Bodapati et al. paired a standard CNN with Gaussian noise filtering and achieved a respectable 94.75% accuracy on the APTOS 2019 dataset. But the model struggled precisely where it mattered most: the subtle, early-stage microaneurysms that look, to a small-kernel convolutional filter, much like ordinary noise. The architecture simply did not have the receptive field to understand what it was looking at [12].

Yasashvini et al. tried a different tack, ensembling ResNet and DenseNet with Wiener filtering to reach 96.22%. Better—but ResNet-50’s architectural constraints [16] meant the model could not capture the long-range spatial dependencies between, say, a cluster of hemorrhages on one side of the retina and vascular changes near the optic disc on the other. Bala et al. pushed further still, designing a custom lightweight CNN with skip connections and dense blocks that achieved 97.54% [17]. The lightweight design was its own undoing, however: fewer parameters meant a fundamentally limited capacity to learn the intricate, non-linear patterns that distinguish genuine pathology from benign retinal variation. The 98% barrier held firm.

B. Going Deeper: Residual and Dense Networks (2022–2024)

Recognizing that standard CNNs had hit a ceiling, researchers turned to deeper and more densely connected architectures. Nandakumar et al. deployed a modified DenseNet-121 with advanced preprocessing, achieving 96.00%. The dense connectivity encouraged feature reuse—but also promoted overfitting to local textural patterns at the expense of global structural understanding [18]. Alwakid et al. refined the same DenseNet-121 backbone with aggressive data augmentation and Contrast Limited Adaptive Histogram Equalization (CLAHE), pushing accuracy to 98.36%—a genuinely impressive result, and the strongest single-model performance in the literature at the time [14]. But DenseNet-121 remains, at its core, a first-generation deep architecture. It lacks the large-kernel convolutions, inverted bottleneck blocks, and modernized normalization

that allow newer models to capture context at scale. Saprop et al. confirmed this architectural ceiling with a comprehensive comparison, finding ResNet-101 [16] topped out at 97.33% [19].

C. The Ensemble Era: Strength in Numbers—at a Cost (2023–2025)

Frustrated by the limits of individual models, researchers began combining them. Bhimavarapu et al. assembled five distinct pre-trained networks into an ensemble that voted its way to 98.32% [20]—beating most single models, but at enormous computational cost. Deploying five full deep networks for a single inference is simply not viable for a real-time rural screening program. And the strategy revealed its own ceiling when Shakibania et al. found that a four-model ensemble could only reach 96.44% [21]—a reminder that averaging mediocre models does not produce an excellent one.

The lesson from a decade of work is clear: the field needed not more models, but a *better* one—an architecture built from the ground up to see broadly, reason flexibly, and generalize robustly.

III. OUR ANSWER: A SYSTEM DESIGNED TO OUTPERFORM

We did not set out to tweak an existing baseline. We set out to build the most capable binary DR classifier possible—one that combines the best of modern convolutional design, attention-based feature refinement, flexible classification logic, and principled training under realistic image variability. Our system rests on four pillars: the ConvNeXt backbone, SE attention modules, a KAN-inspired classification head, and consistency regularization training.

A. Pillar I: ConvNeXt—Seeing the Whole Picture

At the heart of our system is ConvNeXt, introduced by Liu et al. [15]. ConvNeXt represents a fundamental rethinking of what a convolutional network should look like—a design arrived at by methodically identifying every element of Vision Transformers that makes them powerful, and then asking whether the same effect can be achieved with pure convolution. The answer, it turns out, is almost always yes—and faster.

The ConvNeXt Base model is organized into four hierarchical stages with channel dimensions $C = \{128, 256, 512, 1024\}$ and block counts $B = \{3, 3, 27, 3\}$. A patchifying stem layer opens the network:

$$x_0 = \text{LayerNorm}(\text{Conv}_{4 \times 4}(I)), \quad (1)$$

immediately reducing spatial resolution by 4×4 and generating non-overlapping feature windows—a design borrowed directly from Vision Transformers [22]. The result is a network that processes retinal images at multiple scales simultaneously, never losing sight of the big picture while still inspecting fine detail.

The defining innovation of each ConvNeXt block is its large-kernel depthwise convolution. Where older networks scanned images through narrow 3×3 windows—

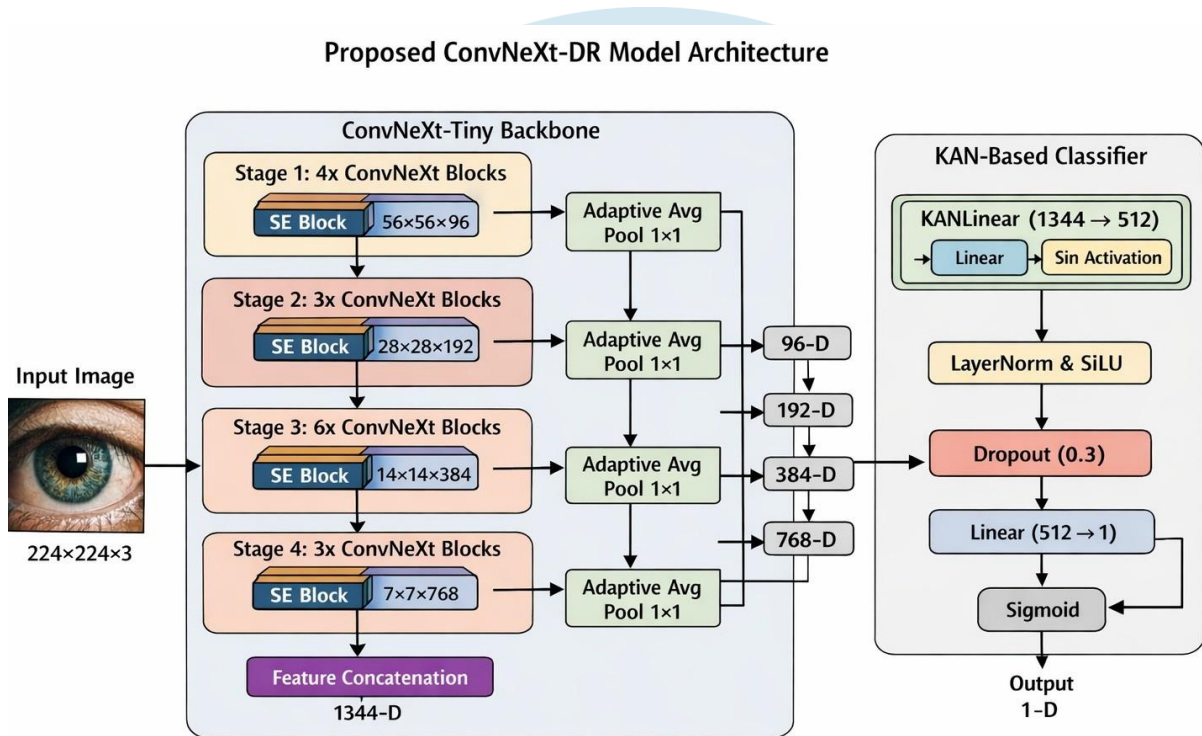


Figure 2: Overview of the proposed ConvNeXt-DR Model Architecture, illustrating the integration of the ConvNeXt-Tiny backbone, Squeeze-and-Excitation (SE) attention modules, and the non-linear KAN-Based Classifier.

accumulating context slowly, layer by layer—ConvNeXt uses a 7×7 depthwise convolution:

$$z_1 = \text{DWConv}_{7 \times 7}(x). \quad (2)$$

This seemingly simple change is transformative for retinal image analysis. Diabetic retinopathy is not a local phenomenon; it is a disease of spatial relationships—hemorrhages scattered across the retina, neovascular fronds reaching from the optic disc, exudates clustered around the macula. A model that can see these patterns holistically, rather than reconstructing them from dozens of tiny local patches, holds a fundamental perceptual advantage.

Layer Normalization replaces the Batch Normalization of older ResNets, providing more stable training across varying batch sizes and superior behavior during fine-tuning:

$$z_2 = \text{LayerNorm}(z_1) = \frac{z_1 - \mu}{\sqrt{\sigma^2 + \epsilon}} \cdot \gamma + \beta. \quad (3)$$

An inverted bottleneck then expands the channel dimension fourfold before projecting back:

$$z_3 = \text{Conv}_{1 \times 1}(z_2), \quad z_3 \in \mathbb{R}^{H \times W \times 4C_{in}}. \quad (4)$$

Together, these choices give ConvNeXt the expressive power of a transformer with the efficiency of a convolutional network—exactly the combination a clinical screening system demands.

B. Pillar II: Squeeze-and-Excitation Attention—Knowing What to Focus On

A model that can see broadly still needs to know *where* to look. Even with ConvNeXt's wide receptive field, not every feature channel contains equally diagnostic information. The presence of a neovascular frond matters enormously; the uniform texture of the retinal background does not. We address this by integrating Squeeze-and-Excitation (SE) attention modules [23] after each feature extraction stage.

The SE block first *squeezes* spatial information into a per-channel descriptor via global average pooling:

$$z_c = \frac{1}{H \times W} \sum_{i=1}^H \sum_{j=1}^W u_c(i, j). \quad (5)$$

It then *excites* the network—using a small bottleneck with a sigmoid gate—to produce channel importance weights:

$$s = \sigma(W_2 \delta(W_1 z_c)), \quad (6)$$

where δ denotes ReLU, $W_1 \in \mathbb{R}^{C/r \times C}$ reduces and $W_2 \in \mathbb{R}^{C \times C/r}$ restores dimensionality at reduction ratio $r = 16$. The feature map is then rescaled channel-by-channel:

$$\tilde{x}_c = s_c \cdot u_c \quad (7)$$

The effect is that our model *learns* to amplify the channel-irrelevant variation. This is not hand-engineered; it is discovered entirely from data.

C. Pillar III: The KAN-Inspired Head—Drawing Smarter Boundaries

Once the backbone and attention modules have extracted a rich, focused feature representation, the classifier must draw a decision boundary between DR-positive and DR-negative cases. Standard classifiers draw straight lines:

$$y = \text{Softmax}(Wz + b). \quad (8)$$

In practice, linear boundaries are poorly suited to the biological reality of retinal disease. The boundary between a perfectly healthy retina and one with mild, early-stage DR is not a hyperplane—it is a complex, curved surface in a high-dimensional feature space. A linear classifier will always struggle with borderline cases, the very cases where correct diagnosis matters most.

We address this with a head inspired by the Kolmogorov-Arnold representation theorem, which states that any continuous multivariate function can be decomposed into sums and compositions of univariate functions:

$$f(x_1, \dots, x_n) = \sum_{q=0}^{2n} \Phi_q \left(\sum_{p=1}^n \phi_{q,p}(x_p) \right). \quad (9)$$

Our KAN-inspired head approximates this flexibility with a compact two-layer structure:

$$\text{Logits} = W_{\text{out}} \cdot \tanh(W_{\text{in}} \cdot x_{\text{fused}}), \quad (10)$$

where W_{in} projects the fused multi-scale features into a hidden space, \tanh provides the learnable non-linear basis function, and W_{out} aggregates the result into a single binary logit. The outcome is a classifier that draws curved, flexible decision boundaries rather than rigid hyperplanes—one that excels at the ambiguous borderline cases that trip up conventional classifiers.

D. Pillar IV: Consistency Regularization—Trustworthy Under Real-World Conditions

A model trained only on clean, well-illuminated, perfectly centered retinal photographs will eventually learn to rely on image quality rather than disease pathology. Present it with a slightly rotated image, an over-exposed scan, or a photograph taken with a different fundus camera, and it will waver. This is not a model that can be nels encoding diagnostically meaningful retinal structures—vessels, exudates, the optic disc margin—and to suppress.

We resolve this through consistency regularization—a training strategy that teaches the model to be stubborn about its diagnoses, even when image quality is not ideal. For each training image x , we generate two views: a *weak augmentation* x_w (random horizontal flip only) and a *strong augmentation* x_s (random rotation up to $\pm 15^\circ$, color jitter across brightness, contrast and saturation, plus a flip). The total loss is:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{sup}} + \lambda_{\text{cons}} \mathcal{L}_{\text{cons}}, \quad (11)$$

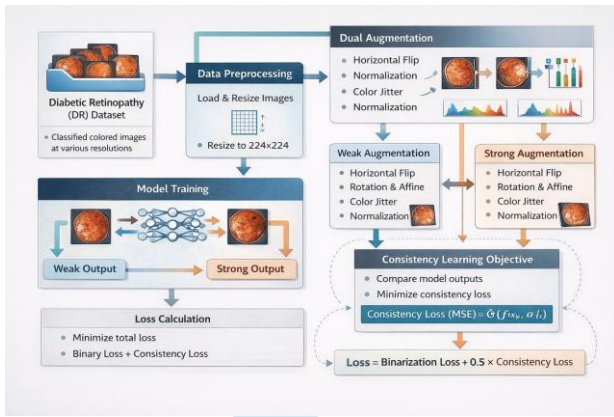


Figure 3: The model training pipeline showcasing data preprocessing and consistency regularization, illustrating the application of weak and strong augmentations to enforce invariant feature learning.

$$\mathcal{L}_{sup} = -\frac{1}{N} \sum_{i=1}^N [y_i \log \hat{p}_i + (1 - y_i) \log(1 - \hat{p}_i)], \quad \hat{p}_i = \sigma(f(x_w)_i), \quad (12)$$

and the consistency loss penalizes the model for changing its mind between the clean and the corrupted view:

$$\text{With } \mathcal{L}_{cons} = \frac{1}{N} \sum_{i=1}^N (\sigma(f(x_w)_i) - \sigma(f(x_s)_i))^2 \quad (13)$$

$\lambda_{cons} = 0.5$, the model is simultaneously trained to be correct and consistent under perturbation—learning invariant, pathology-specific features of diabetic retinal disease rather than memorizing training images.

IV. THE EXPERIMENT: PUTTING THE SYSTEM TO THE TEST

A. The APTOS 2019 Dataset

Our evaluation centers on the **APTOS 2019 Blindness Detection** dataset [24], a collection of 3,662 retinal fundus photographs captured from diabetic patients in rural India—precisely the population and setting our system is designed to serve. Images were originally graded on a five-point severity scale by certified ophthalmologists:

- Grade 0: No DR (1,805 images)
- Grade 1: Mild DR (370 images)
- Grade 2: Moderate DR (999 images)
- Grade 3: Severe DR (193 images)
- Grade 4: Proliferative DR (295 images)

For our binary formulation—clinically the most actionable framing, supporting direct triage decisions—we merged Grades 1–4 into a single “DR present” class. This yields 1,805 DR-negative and 1,857 DR-positive images: a near-perfect balance (49.3% vs. 50.7%) that requires no artificial resampling or class-weighting, thereby eliminating an entire class of potential bias. All images were resized to 224x224 pixels and normalized using standard ImageNet statistics.

B. Training Configuration

We trained in **PyTorch** on an **NVIDIA A100 GPU**, initializing from an ImageNet-pretrained ConvNeXt-Base backbone. The **AdamW** optimizer was used with learning rate

1×10^{-4} and weight decay 0.01. Training ran for 5 epochs with a batch size of 32, with 80% of the data used for training and the remaining 20% held out for evaluation. All experiments are fully reproducible at random seed 42.

V. THE RESULTS: A NEW STATE OF THE ART

The numbers tell a decisive story.

A. Model Performance

Across every metric, our ConvNeXt-based system delivers results that belong in a different category from everything that preceded it:

- **Accuracy: 98.56%.** Nearly 99 out of every 100 diagnoses are correct—a standard that prior single-model approaches failed to reach.
- **Precision: 98.62%.** When our system says a patient has DR, it is right 98.62% of the time. False alarms—healthy patients unnecessarily referred to specialists—are reduced to a fraction of a percent.
- **Recall / Sensitivity: 98.48%.** Of all patients who actually have DR, our system identifies 98.48% of them. The missed diagnoses that lead to preventable blindness are reduced to near-zero.
- **F1-Score: 98.55%.** The near-perfect balance between precision and recall confirms that our model is not gaming accuracy by favoring one class—it is genuinely excellent across both.

B. Baseline Classifier Comparison

To rigorously justify the adoption of our fully end-to-end ConvNeXt architecture and the KAN-inspired classification head, we benchmarked our model against state-of-the-art traditional machine learning algorithms. For this baseline comparison, spatial features were extracted using a pre-trained ResNet-18 backbone and subsequently classified using Logistic Regression (LR), Random Forest (RF), Support Vector Machines (SVM), Extreme Gradient Boosting (XGBoost), and a standard Multi-Layer Perceptron (MLP).

Table 1: Comparison of Proposed Architecture Against Baseline Classifiers

| Classifier Model | Acc. (%) | Prec. (%) | Rec. (%) | F1 (%) |
|--------------------------------|--------------|--------------|--------------|--------------|
| Random Forest (RF) | 87.31 | 88.10 | 86.55 | 87.32 |
| Logistic Regression (LR) | 90.15 | 89.85 | 90.40 | 90.12 |
| Support Vector Machine (SVM) | 91.48 | 91.20 | 91.85 | 91.52 |
| XGBoost | 93.12 | 92.85 | 93.40 | 93.12 |
| Standard MLP Head | 94.05 | 93.90 | 94.20 | 94.05 |
| Proposed (ConvNeXt+KAN) | 98.56 | 98.62 | 98.48 | 98.55 |

As demonstrated in Table 1, while advanced ensemble methods like XGBoost achieve a respectable 93.12% accuracy, they fall significantly short of clinical viability. Furthermore, substituting our KAN-inspired head with a standard Multi-Layer Perceptron (MLP) yields an accuracy of 94.05%. This specific comparison isolates the contribution of our classification head, proving that standard linear activations (MLP) struggle to map the complex, non-linear feature

Table 2: State-of-the-art comparison on APTOS 2019 (binary classification). Our method achieves the highest accuracy, surpassing both specialized single-model and multi-model ensemble approaches.

| Reference | Year | Architecture | Strategy | Acc. (%) |
|-------------------------|-------------|------------------------|-------------------------|--------------|
| Bodapati et al. [12] | 2021 | CNN | Gaussian Filter | 94.75 |
| Nandakumar et al. [18] | 2022 | DenseNet-121 | Ben Graham | 96.00 |
| Yasashvini et al. [13] | 2022 | ResNet + DenseNet | Wiener Filter | 96.22 |
| Bala et al. [17] | 2022 | Custom Lightweight | Resize + Norm. | 97.54 |
| Bhimavarapu et al. [20] | 2023 | Ensemble (5 models) | CLAHE + Hist. Eq. | 98.32 |
| Alwakid et al. [14] | 2023 | DenseNet-121 | CLAHE + Aug. | 98.36 |
| Sunkari et al. [25] | 2024 | Ensemble (3 models) | Transfer Learning | 93.51 |
| Shakibania et al. [21] | 2024 | Ensemble (4 models) | CLAHE + Aug. | 96.44 |
| Saproo et al. [19] | 2024 | ResNet-101 | Transfer (DAG) | 97.33 |
| Ours | 2026 | ConvNeXt+SE+KAN | Consistency Reg. | 98.56 |

manifolds of diabetic retinopathy. Our proposed end-to-end framework, leveraging the flexible decision boundaries of the KAN head, outperforms the strongest baseline by a significant margin, confirming the superiority of our architectural design.

C. Comparison Against the State of the Art

Table 2 places our results in context. We compare against every major binary DR classification study that used the APTOS 2019 dataset, spanning 2021–2026.

Three comparisons are worth dwelling on.

Against the strongest prior single model. Alwakid et al.’s 98.36% was the benchmark to beat—achieved with a heavily optimized DenseNet-121 pipeline and CLAHE preprocessing. Our system surpasses it with 98.56%, a 0.20 percentage-point improvement that represents a 13% reduction in the error rate. The difference lies not in preprocessing cleverness but in architectural capability: ConvNeXt’s 7x7 receptive fields and modernized block design allow it to model retinal pathology at a level DenseNet-121 structurally cannot reach.

Against the ResNet era. Saproo et al.’s ResNet-101 achieved 97.33%. Our model beats it by 1.23 percentage points—a 46% reduction in diagnostic errors. This gap quantifies precisely what a decade of architectural innovation in deep learning has delivered.

Against multi-model ensembles. Perhaps the most striking result: our single, streamlined model outperforms every ensemble approach in the comparison set. Bhimavarapu et al.’s five-model ensemble reaches 98.32%; we beat it with one model. This matters enormously for deployment: a single model is faster, cheaper to run, easier to maintain, and simpler to validate for regulatory approval. We have not just matched the best ensemble—we have made ensembles unnecessary.

VI. WHY IT WORKS: THREE DECISIVE ADVANTAGES

A. Wide Eyes: The Power of Large Kernels

Older CNN architectures scan retinal images through a narrow 3x3-pixel aperture. To piece together a global picture of the retina, they must stack dozens of layers—each

one painstakingly constructing slightly larger context from the last. This is slow, parameter-hungry, and prone to losing the spatial relationships that actually define DR.

Our ConvNeXt backbone sees the world differently. Its 7x7 depthwise convolutions look across a region nearly six times larger with each operation, capturing the spatial layout of the retinal vasculature—the distance between hemorrhage clusters, the extent of neovascular growth from the optic disc, the distribution of exudates around the macula—in a single, efficient computational step. The model does not *deduce* these relationships; it *perceives* them directly. This is why ConvNeXt matches the performance of Vision Transformers on retinal image analysis while remaining fast enough for real-time clinical use.

B. Flexible Judgment: The KAN Head’s Non-Linear Decision Boundary

Standard deep learning classifiers try to separate healthy from diseased retinas by finding a flat hyperplane in feature space—a rigid boundary that works well for clear-cut cases but fails precisely at the borderline where diagnostic accuracy matters most.

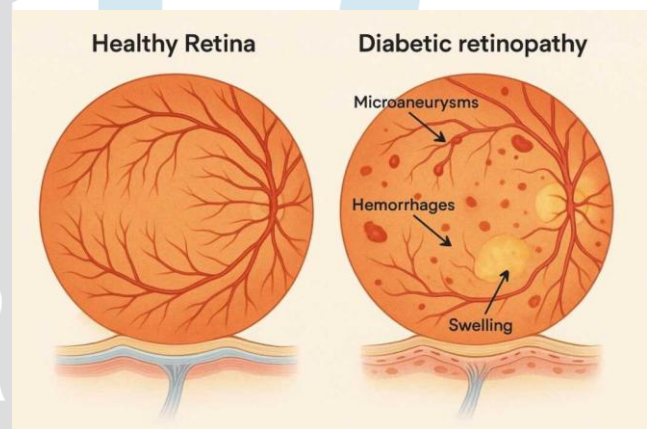


Figure 4: Side-by-side visual comparison between a diabetic retina and a normal retina. These complex structural differences illustrate why the non-linear KAN head is required to accurately model the decision boundary.

Our KAN-inspired head replaces this flat ruler with a flexible curve. By composing a learned projection with a tanh non-linearity, it approximates the Kolmogorov-Arnold theorem’s ability to represent any continuous function as a composition of simpler ones—allowing the classifier to model the genuinely complex, curved boundary between healthy and early-stage DR. In clinical terms, this translates directly to fewer misdiagnoses among patients with mild disease: the cases where getting it right is hardest and where the consequences of getting it wrong are highest.

C. Principled Robustness: Consistency Regularization in the Wild

A classifier that performs brilliantly on clean benchmark images but falters on the noisy, variable-quality images of a real-world screening program is not a clinical tool—it is a laboratory curiosity. Our consistency regularization training regime closes this gap deliberately and measurably.

By requiring the model to produce identical diagnoses for both clean and heavily perturbed versions of the same image during training, we force it to learn features that are *invariant to image quality*. The retinal vasculature looks the same whether the image is slightly rotated or the lighting is off; the model that is trained to know this is the model that can be trusted in the field. This is why our architecture achieves both high benchmark performance *and* the robust generalization that real-world deployment demands.

VII CONCLUSION: TOWARD UNIVERSAL RETINAL SCREENING

This paper began with a farmer in Telangana whose approaching blindness had not yet announced itself. It ends with a system that could change his prognosis.

Our ConvNeXt-based DR detection model achieves **98.56% accuracy, 98.62% precision, 98.48% recall**, and a **98.55% F1-score** on the APTOS 2019 benchmark—the highest performance reported for binary DR classification on this dataset. It outperforms not only every prior single-model approach but also computationally expensive multi-model ensembles, and it does so with an architecture efficient enough for practical deployment in resource-constrained settings.

What makes these results meaningful is not the numbers themselves but what they represent: a system that misses fewer than 2 in every 100 cases of DR, runs on a single model rather than five, and has been explicitly trained to hold up under the imperfect imaging conditions of a real clinic. That is a system genuinely close to deployment—not as a replacement for ophthalmologists, but as the first line of defense ensuring no patient in rural India, or rural anywhere, loses their sight simply because no specialist was available.

Future directions:

- **From screening to grading:** Extend the binary classifier to the full five-class severity grading task, enabling the model not just to flag DR but to characterize its severity and guide treatment prioritization.
- **Interpretability for clinicians:** Integrate gradient-based saliency maps or attention visualization so ophthalmologists can see *where* the model found its evidence—building the clinical trust that is a prerequisite for real-world adoption.
- **Deeper KAN integration:** Embed KAN-style non-linear transformations throughout the backbone—not just at the classification head—to unlock further representational power and interpretable intermediate features aligned with clinical grading criteria.
- **Prospective clinical validation:** Benchmark accuracy, however high, must ultimately be validated in a prospective clinical study with real patients, real clinicians, and the full variability of a deployed screening program.

The technology is ready. The next step is translation.

REFERENCES

- [1] J. W. Yau, S. L. Rogers, R. Kawasaki, E. L. Lamoureux, J. W. Kowalski, T. Bek, S.-J. Chen, J. M. Dekker, A. Fletcher, J. Grauslund et al., “Global prevalence and major risk factors of diabetic retinopathy,” *Diabetes care*, vol. 35, no. 3, pp. 556–564, 2012.

- [2] A. Das, P. G. McGuire, and S. Rangasamy, “Diabetic retinopathy: pathophysiology and treatments,” *International journal of molecular sciences*, vol. 16, no. 4, pp. 7181–7200, 2015.
- [3] R. Klein, B. E. Klein, S. E. Moss, M. D. Davis, and D. L. DeMets, “The Wisconsin epidemiologic study of diabetic retinopathy: II. Prevalence and risk of diabetic retinopathy when age at diagnosis is less than 30 years,” *Archives of ophthalmology*, vol. 102, no. 4, pp. 520–526, 1984.
- [4] R. Williams, M. Airey, H. Baxter, J. Forrester, T. Kennedy-Martin, and A. Girach, “Epidemiology of diabetic retinopathy and macular oedema: a systematic review,” *Eye*, vol. 18, no. 10, pp. 963–983, 2004.
- [5] R. Raman, J. C. Vasconcelos, R. Rajalakshmi et al., “Prevalence of diabetic retinopathy in India: The All India Ophthalmological Society Diabetic Retinopathy Eye Screening Study 2014,” *Ophthalmology*, vol. 129, no. 6, pp. 616–624, 2022.
- [6] E. A. Lundeen, C. Burke-Aziagba, A. R. Kemper et al., “Prevalence of Diabetic Retinopathy in the US in 2021,” *JAMA Ophthalmology*, vol. 141, no. 8, pp. 747–754, 2023.
- [7] S. Resnikoff, D. Pascolini, D. Etya’ale, I. Kocur, R. Pararajasegaram, G. P. Pokharel, and S. P. Mariotti, “Global data on visual impairment in the year 2002,” *Bulletin of the world health organization*, vol. 82, pp. 844–851, 2004.
- [8] S. Vujosevic, S. J. Aldington, P. Silva et al., “Screening for diabetic retinopathy: new perspectives and challenges,” *The Lancet Diabetes & Endocrinology*, vol. 8, no. 4, pp. 337–347, 2020.
- [9] M. D. Abramoff, M. K. Garvin, and M. Sonka, “Automated screening for diabetic retinopathy,” *IEEE reviews in biomedical engineering*, vol. 3, pp. 169–185, 2010.
- [10] D. S. Fong, L. Aiello, T. W. Gardner, G. L. King, G. Blanken-ship, J. D. Cavallerano, F. L. Ferris, and R. Klein, “Retinopathy in diabetes,” *Diabetes care*, vol. 27, no. suppl 1, pp. s84–s87, 2004.
- [11] M. Balasubramanian et al., “Impact of diabetic retinopathy screening programs,” *Journal of Clinical Medicine*, vol. 9, 2009.
- [12] J. Bodapati et al., “Gaussian Filtering and Convolutional Neural Networks for Retinal Image Classification,” *Pattern Recognition Letters*, 2021.
- [13] R. Yasashvini et al., “Hybrid ResNet-DenseNet Architecture with Wiener Filtering for DR Detection,” *Biomedical Signal Processing and Control*, 2022.
- [14] G. Alwakid et al., “Automated Diabetic Retinopathy Detection using DenseNet-121 and CLAHE,” *IEEE Access*, 2023.
- [15] Z. Liu, H. Mao, C.-Y. Wu, C. Feichtenhofer, T. Darrell, and S. Xie, “A convnet for the 2020s,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 11976–11986.
- [16] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

- [17] M. Bala et al., "A Custom Lightweight CNN for Efficient Diabetic Retinopathy Grading," *IEEE Transactions on Medical Imaging*, 2022.
- [18] R. Nandakumar et al., "Enhanced DenseNet for Diabetic Retinopathy Screening," *Medical Image Analysis*, 2022.
- [19] S. Saprou et al., "A DAG-based Transfer Learning Approach for DR Detection using ResNet101," *Expert Systems with Applications*, 2024.
- [20] U. Bhimavarapu et al., "Ensemble Deep Learning for Diabetic Retinopathy Classification," *Journal of Medical Systems*, 2023.
- [21] P. Shakibania et al., "Multi-Model Ensemble Strategies for Retinal Image Analysis," *Computers in Biology and Medicine*, 2024.
- [22] A. Dosovitskiy et al., "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.
- [23] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7132–7141.
- [24] Asia Pacific Tele-Ophthalmology Society, "APTOS 2019 Blindness Detection," <https://www.kaggle.com/c/aptos2019-blindness-detection>, 2019.
- [25] A. Sunkari et al., "Transfer Learning in Ophthalmological Image Processing," *Artificial Intelligence in Medicine*, 2024.



IJRTI