# NEURO-GUARD: Neuro-Symbolic Generalization and Unbiased Adaptive Routing for Diagnostics - Explainable Medical AI

Midhat Urooj
Arizona State University
Tempe, AZ, USA
murooj@asu.edu

Ayan Banerjee
Arizona State University
Tempe, AZ, USA
abanerj3@asu.edu

Sandeep Gupta
Arizona State University
Tempe, AZ, USA
Sandeep.Gupta@asu.edu

## Abstract

*Accurate yet interpretable image-based diagnosis remains a central challenge in medical AI, particularly in settings with limited data, subtle visual patterns, and high-stakes clinical decisions. However, most current vision models produce black-box predictions with limited generalizability and poor real-world usability. We present **NEURO-GUARD**, a novel framework that combines Vision Transformers (ViTs) with knowledge-guided reasoning to enhance **performance, transparency**, and **cross-domain generalization**. NEURO-GUARD incorporates a **retrieval-augmented generation (RAG)** mechanism for language-driven self-verification, in which a large language model (LLM) iteratively generates, evaluates, and refines feature extraction code for medical images. By leveraging clinical guidelines and expert knowledge, this LLM-guided module progressively improves feature detection and classification, outperforming purely data-driven baselines. Extensive evaluations on diabetic retinopathy classification across four benchmark datasets (APTOS, EyePACS, Messidor-1, Messidor-2) show that NEURO-GUARD improves accuracy by **6.2%** over a ViT-only model (84.69% vs. 78.4% [3]) and achieves a **5%** gain in domain generalization. Further experiments on MRI-based seizure detection confirm its cross-domain robustness, consistently surpassing existing baselines. Notably, NEURO-GUARD bridges the gap between symbolic medical reasoning and subsymbolic feature learning, demonstrating robust generalization across multiple datasets while achieving **state-of-the-art performance**.*

## 1. Introduction

Medical imaging plays a crucial role in disease diagnosis and treatment planning, particularly in conditions such as diabetic retinopathy (DR), tumor detection, and neurodegenerative disorders. Recent advances in deep learning,
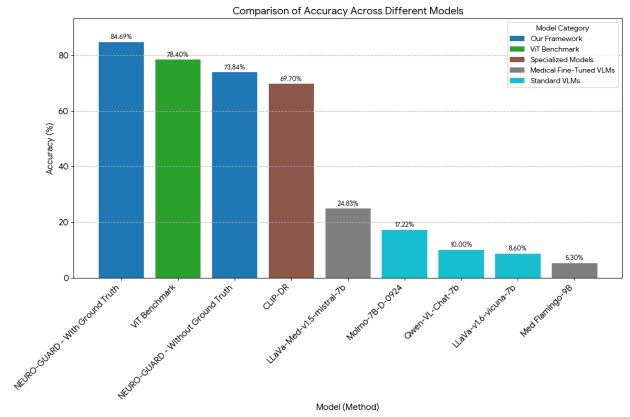


Figure 1. Performance comparison of existing models versus the NEURO-GUARD framework for 5-stage Diabetic Retinopathy classification.

particularly Vision Transformers (ViTs) and Convolutional Neural Networks (CNNs), have significantly improved diagnostic accuracy [3, 14]. However, their black-box nature limits clinical adoption due to a lack of interpretability, making it challenging for clinicians to validate AI-driven decisions. Additionally, these models suffer from domain shift vulnerabilities, struggling to generalize across imaging datasets with diverse acquisition protocols and patient demographics [22, 26]. Given these challenges, an ideal medical AI framework should not only provide high accuracy but also generate clinically interpretable decisions by integrating structured domain knowledge into its reasoning process.

Existing explainability techniques, such as Gradient-weighted Class Activation Mapping (Grad-CAM) [13] and Shapley Additive Explanations (SHAP) [8], provide post-hoc feature attribution but remain static, heuristic-based, and disconnected from the model's decision logic. Hybrid approaches incorporating attention mechanisms and uncertainty estimation attempt to improve interpretabil-
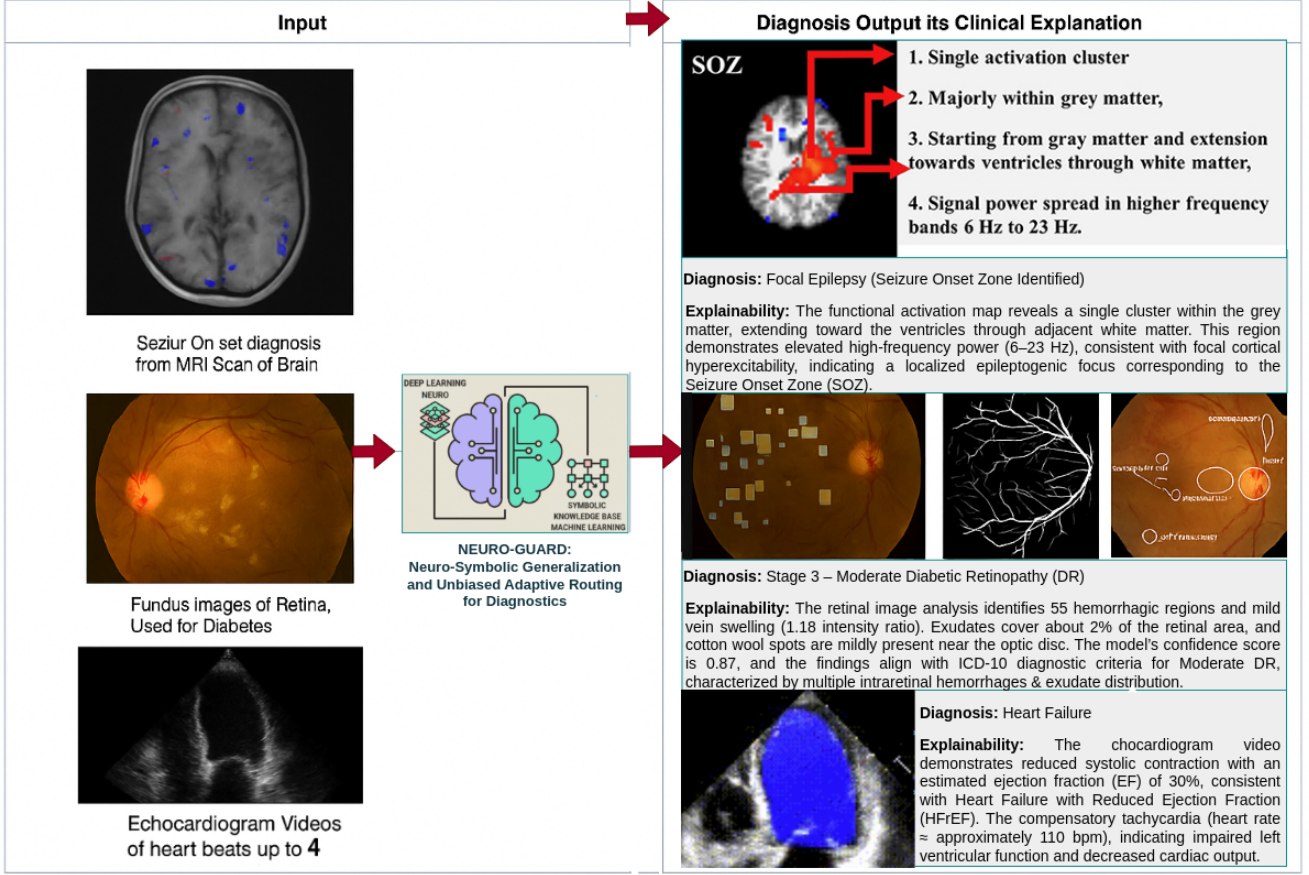
Figure 2. Overview of the NEURO-GUARD framework. The system integrates medical knowledge with multimodal imaging to enhance disease classification and provide clinically aligned, interpretable explanations with spatial localization.

ity [18, 20], but they fail to integrate structured medical knowledge, limiting their ability to generalize across datasets. Reinforcement learning (RL) and meta-learning frameworks [10] enable adaptive learning, yet they lack mechanisms to ground AI decisions in clinical reasoning, reducing their reliability in real-world medical applications.

To address these limitations, we propose NEURO-GUARD, a novel framework that fuses language-grounded reasoning with state-of-the-art visual recognition to enable intrinsically interpretable medical image diagnosis. In contrast to prior systems that only add interpretability after the fact, NEURO-GUARD tightly integrates a clinical knowledge base and reasoning module into the model's inference pipeline. This is achieved through a modular architecture combining a self-supervised ViT-based image encoder with a knowledge-guided language model that jointly analyzes images and textual information. Crucially, our approach leverages retrieval-augmented generation to dynamically draw on external biomedical sources (e.g., literature, guidelines) for case-specific knowledge, and uses an LLM-based code synthesis engine to translate this knowl-

edge into executable image analysis steps. A prompt-driven self-verification loop, optimized via reinforcement learning, compels the model to iteratively check and refine its outputs greatly reducing hallucinations and aligning final predictions with clinical guidelines. Through this design, NEURO-GUARD shifts interpretability from a post-hoc exercise to an intrinsic property of the model's predictions, as the reasoning is conducted in natural language and grounded in real clinical criteria from the start. In essence, our framework bridges the gap between symbolic medical knowledge and subsymbolic vision features, allowing the model to explain why and how it arrives at a diagnosis in terms familiar to human experts as shown in Figure 4.

## 1.1. Contributions

We design a **hybrid inference pipeline** that fuses deep-learning predictions with **knowledge-driven classifiers**, enabling transparency while maintaining high diagnostic performance. In contrast to existing VLMs and LLMs which often produce confident yet inaccurate and hallucinated explanations, as demonstrated in our Phase 1 and
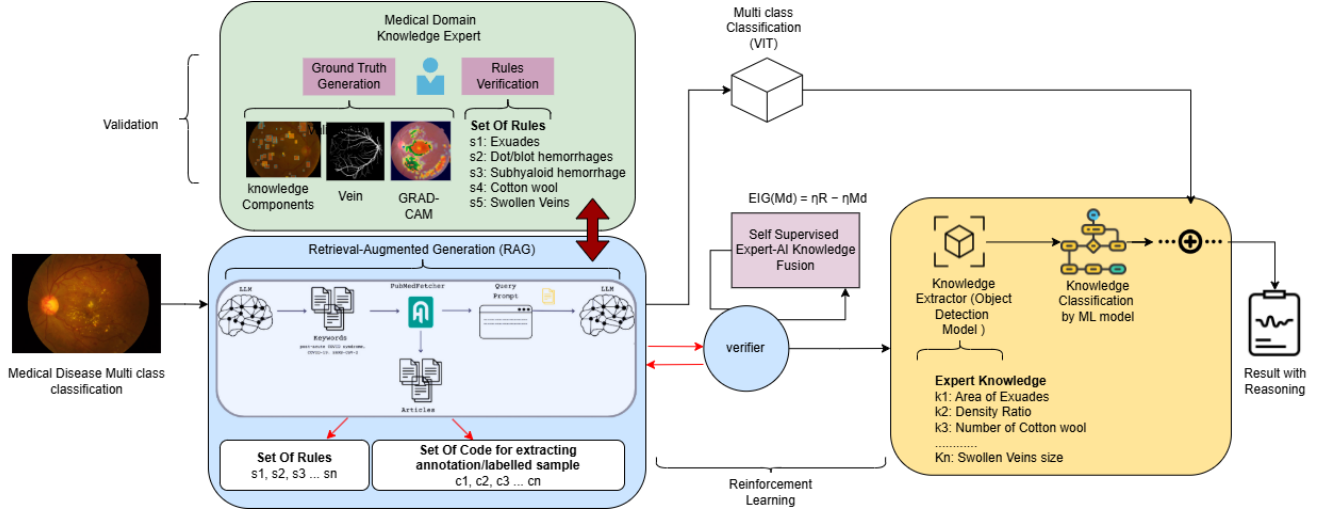
Figure 3. **NEURO-GUARD Framework for Knowledge-Driven Medical AI.** The NEURO-GUARD pipeline integrates RAG-based knowledge extraction, reinforcement learning-based self-verification, and multi-class classification.

Phase 2 experiments. NEURO-GUARD grounds every decision in structured medical knowledge. We introduce the **first medical-imaging pipeline that transforms clinical guidelines and expert rules into executable, verifiable code**, integrating this symbolic reasoning directly with **Vision Transformer (ViT) feature learning**. Using a **multi-stage Retrieval-Augmented Generation (RAG) process** grounded in peer-reviewed medical literature and disease-specific protocols, our system constructs dynamic rule bases that guide pixel-level lesion detection with stronger clinical consistency. To further ensure reliability, we develop an **entropy-based reinforcement learning self-verification loop** that iteratively refines code-generated feature extractors, substantially reducing hallucinations and improving localization accuracy. Through comprehensive evaluations on **diabetic retinopathy datasets** (APTOS, EyePACS) as well as **MRI-based seizure detection**, we demonstrate that NEURO-GUARD achieves robust generalizability across diseases, imaging modalities, and clinical conditions, establishing a new direction for interpretable and trustworthy medical AI.

## 2. Related Work

### 2.1. Interpretable Models vs. Black-Box Paradigms

The trade-off between predictive performance and model transparency remains a defining challenge in medical artificial intelligence. While deep learning architectures, such as Convolutional Neural Networks (CNNs) [14] and Vision Transformers (ViTs) [3], have achieved state-of-the-art performance in diagnostic tasks, they operate as opaque "black boxes." Early attempts to mitigate this opacity relied on post-hoc explanation methods like Grad-CAM [13], which

highlight salient image regions. However, these methods have been criticized for providing fragile approximations that often fail to capture the model's true decision logic [12]. Rudin [12] argues that for high-stakes decisions, reliance on post-hoc explanations is insufficient and advocates for inherently interpretable architectures.

In response, Concept Bottleneck Models (CBMs) [7] were proposed to align latent representations with human-understandable concepts (e.g., "bone spur" or "effusion") prior to classification. While CBMs offer intrinsic interpretability, they typically require dense, costly concept annotations and often suffer from a performance gap compared to end-to-end models [24]. Alternative strategies employing reinforcement learning [22] offer dynamic adaptation but often lack mechanisms to ground learned representations in explicit domain knowledge. Unlike these static or purely data-driven approaches, our framework leverages the dynamic reasoning capabilities of Large Language Models (LLMs) to generate interpretable feature extractors without requiring pixel-level concept supervision.

### 2.2. Vision-Language Models in Medicine

The integration of linguistic knowledge into medical imaging has been accelerated by Vision-Language Models (VLMs). Building on the contrastive learning paradigm of CLIP [11], domain-specific models such as MedCLIP [21] and Med-PaLM 2 [15] have demonstrated remarkable capabilities in zero-shot classification and report generation. These models encode vast medical ontologies, enabling them to process complex diagnostic queries.

However, significant limitations persist. Recent studies indicate that generative medical agents frequently exhibit "hallucinations," producing plausible but factually incor-

rect findings due to a lack of grounding in pixel-level evidence [2]. Furthermore, most VLMs process information through separate, static encoders, losing the fine-grained feature alignment necessary for verifying subtle biomarkers [25]. NEURO-GUARD addresses this by moving beyond static embeddings, using an agentic framework to actively query and verify visual features.

## 2.3. Neuro-Symbolic and Agentic Frameworks

To bridge the gap between symbolic reasoning and sub-symbolic perception, recent works in general computer vision have proposed "code-as-policy" approaches. Systems such as VisProg [5] and ViperGPT [16] utilize LLMs to decompose complex visual queries into executable Python programs, invoking vision primitives to solve tasks without specific training. These frameworks demonstrate that LLMs can act as reasoning engines to orchestrate vision modules effectively.

Despite their success in general domains, these agents lack the specialized clinical knowledge required for medical diagnostics. General-purpose code generation often fails to synthesize the precise, domain-specific subroutines needed to detect pathological features (e.g., distinguishing "microaneurysms" from "hemorrhages"). Our work extends the neuro-symbolic paradigm by integrating Retrieval-Augmented Generation (RAG) with code generation, synthesizing expert clinical guidelines into executable logic that is both performant and verifiable.

## 3. Proposed method

Our proposed pipeline, NEURO-GUARD, autonomously identifies and refines domain-specific features by integrating symbolic medical knowledge with subsymbolic learning, ensuring strong alignment with clinical standards. As depicted in Figure 3, NEURO-GUARD employs a Retrieval-Augmented Generation (RAG) mechanism to extract structured medical knowledge from sources such as PubMed and clinical guidelines. For instance, in the context of diabetic retinopathy (DR), the system retrieves rules indicating that hemorrhages, exudates, and swollen veins are key visual markers. This retrieved knowledge is passed to a large language model (LLM) in a multi-prompt sequence. In Prompt 1, the LLM consolidates disease-specific information into a structured clinical rule base, detailing relevant features and diagnostic criteria. In Prompt 2, the LLM utilizes this rule base to generate executable Python code for feature detection, embedding reinforcement learning (RL) parameters to guide initial predictions. In Prompt 3, performance feedback including metrics such as Intersection-over-Union (IoU), precision, and recall is used to iteratively refine the generated code via RL, enhancing alignment with human annotations or model confidence. A self-verification module then evaluates whether the extracted
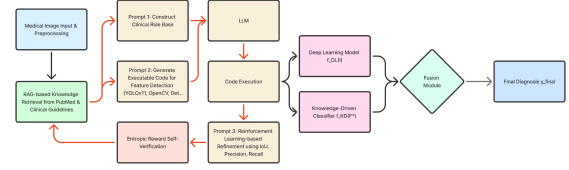


Figure 4. Flow diagram of the NEURO-GUARD framework

features conform to expected clinical patterns using entropic reward signals. Ultimately, NEURO-GUARD produces interpretable and generalizable outputs by harmonizing LLM-driven symbolic reasoning with deep learning-based feature recognition.

### 3.1. Mathematical Formulation

Given an input retinal image $I$, we first apply a feature detection module $\mathcal{F}$ to extract a comprehensive set of clinically relevant features $F = \{f_1, f_2, \ldots, f_n\}$. Each feature $f_i$ corresponds to distinct anatomical or pathological markers, such as exudates, hemorrhages, and cotton wool spots. The feature set $F$ is mapped into a vectorized representation $v \in \mathbb{R}^d$ through a domain-specific embedding function $\varphi : F \to \mathbb{R}^d$.

Simultaneously, we retrieve and structure clinical knowledge rules $R$, which include visual representation rules $R_V = \{r_{v1}, r_{v2}, \ldots, r_{vm}\}$ and demographic knowledge rules $R_D = \{r_{d1}, r_{d2}, \ldots, r_{dl}\}$. These rules are dynamically extracted from authoritative clinical sources such as PubMed using a Retrieval-Augmented Generation (RAG) framework. Each visual rule $r_{vi}$ is directly associated with a corresponding subset of image-derived features $F_{rvi}$, and demographic rules $r_{dj}$ are linked to patient-specific metadata features $F_{rdj}$, forming an integrated knowledge extraction set:

$$F = \{F_{rv1}, F_{rv2}, \ldots, F_{rvm}, F_{rd1}, F_{rd2}, \ldots, F_{rdl}\}$$

### 3.2. Visual Feature Extraction and Self-Verification

Large Language Models (LLMs) are prompted to generate executable visual feature extraction codes $C = \{c_1, c_2, \ldots, c_m\}$, each specifically tailored to detect features corresponding to visual rules $r_{vi}$. These codes are instantiated through vision frameworks such as OpenCV, YOLOv11, or Detectron2. Extracted visual features $X_i$ from an image $I$ are quantitatively validated against ground truth annotations $Y_i$ using Intersection-over-Union (IoU) metrics:

$$\text{IoU}(X_i, Y_i) = \frac{|X_i \cap Y_i|}{|X_i \cup Y_i|}$$

To enhance the reliability and accuracy of feature extraction, we introduce an entropic reward-based self-

verification mechanism. The entropic gain $E_i$ quantifies the reduction in uncertainty associated with feature extraction:

$$E_i = -\sum_{k=1}^{K} p_k \log p_k$$

where $p_k$ denotes the empirical probability of correct extraction across $K$ randomly selected validation images. If the entropy $E_i$ is below a predefined threshold $\tau$, iterative refinement of the extraction code $c_i$ is performed using reinforcement learning (RL)-guided prompt tuning:

$$c_i^{(t+1)} = \text{LLM}\left(E_i, c_i^{(t)}\right), \quad \text{until} \quad \max_t E_i^{(t)} \geq \tau$$

Upon convergence, verified visual extraction codes $C^*$ yield the optimized visual knowledge set $F_{RV}^*$. Together with the demographic knowledge set $F_{RD}$, they form the refined feature set:

$$F^* = \{F_{RV}^*, F_{RD}\}$$

### 3.3. Classification and Decision Fusion

The refined feature set $F^*$ is input into a knowledge-driven classification model $f_{KD}$, producing a predicted disease stage $y_{KD} = f_{KD}(F^*)$. Concurrently, a deep learning-based model $f_{DL}$, trained on extensive labeled retinal images, processes the input image $I$, generating a prediction $y_{DL} = f_{DL}(I)$ with associated confidence $s_{DL}$.

The final diagnosis $y_{final}$ is obtained by integrating both predictions through a fusion function $f_{fusion}$:

$$y_{final} = \begin{cases} y_{DL}, & \text{if } s_{DL} \geq s_{KD} \\ y_{KD}, & \text{otherwise.} \end{cases}$$

### 3.4. Reinforcement Learning-Based Optimization

To systematically optimize detection accuracy, we utilize an RL-based parameter tuning approach under two distinct feedback environments:

**Supervised Artifact Detection - With Ground Truth Annotations**   In this scenario, RL directly optimizes YOLO confidence thresholds:
- **State Representation**: Discretized YOLO confidence threshold.
- **Action Space**: Adjustments by increments {-0.05, 0, +0.05}.
- **Reward Function**: IoU between predicted and annotated bounding boxes.
- **Q-learning Update**: Standard Q-learning algorithm with parameters $\alpha = 0.1$, $\gamma = 0.9$, determined via empirical sensitivity analysis.

**Unsupervised Knowledge Extraction - Without Ground Truth Annotations**   OpenCV parameters optimization is performed heuristically:
- **State Representation**: Tuple of CLAHE clip limit and optic disc masking threshold.
- **Action Space**: Incremental parameter adjustments.
- **Reward Function**: Heuristic measure based on the accuracy of detected feature counts relative to clinically expected targets.
- **Training and Q-learning Update**: Mirrors the ground truth approach, maintaining consistent hyperparameter values for stability and generalization.

## 4. Experimental Setup and Evaluation

Diabetic Retinopathy (DR) remains a leading cause of visual impairment, necessitating diagnostic methods that are not only accurate but also interpretable. While deep learning has achieved success in classification [9], clinical adoption requires the precise localization of pathological features such as exudates and hemorrhages to justify decisions. Our objective was to achieve this reasoning and explainability. However, our initial experiments with state-of-the-art foundation models revealed critical reliability gaps, motivating the development of the **NEURO-GUARD** framework. The experimental setup shows the full journey of experimental phases from getting results based on hallucinations to highly accurate results with correct Reasoning from our NEURO-GUARD framework.

**Phase 1: Failure of Zero-Shot Vision-Language Models.** We first evaluated off-the-shelf Vision-Language Models (VLMs) on the APTOS and EyePACS datasets. Inspired by the CARES benchmark for trustworthiness [23], we tested zero-shot and few-shot capabilities. As detailed in Table 3, these models struggled significantly with medical granularity:
- **CLIP & MedCLIP:** These models relied on holistic semantics and failed to capture subtle lesions, often misclassifying severe DR as mild due to a lack of object-level grounding [1].
- **Grounding-DINO:** Produced frequent false positives by incorrectly labeling non-pathological artifacts (e.g., optic disc reflections) as lesions [17].
- **InstructBLIP:** Lacked the spatial resolution required to detect small markers like microaneurysms.

These findings align with recent studies on medical hallucinations, where models generate convincing but factually incorrect diagnostic captions [4, 6].

**Phase 2: Limitations of Direct LLM-Based Feature Extraction.**   Given the localization failures of VLMs, we explored using Large Language Models (LLMs) like Gemini,

Table 1. Performance comparison of classification models on the Aptos and Eye Pacs datasets. Models are tested with supervised knowledge and without supervised knowledge to assess their interpretability impact. The best-performing values per row are **bolded**, and the highest-performing model per dataset is highlighted.

| Framework | Model | Dataset | Val. Acc. | Test Acc. | Precision | Recall | F1-Score | Weighted Prec. | Weighted F1 |
|---|---|---|---|---|---|---|---|---|---|
| NEURO-GUARD (Unsupervised Knowledge Extraction) | Logistic Regression | Aptos | 0.6019 | 0.6424 | 0.25 | 0.33 | 0.28 | 0.55 | 0.58 |
| | Random Forest | Aptos | 0.7038 | **0.7384** | 0.55 | 0.47 | 0.49 | 0.71 | 0.71 |
| | SVM | Aptos | 0.6083 | 0.6556 | 0.26 | 0.34 | 0.28 | 0.56 | 0.59 |
| | Gradient Boosting | Aptos | 0.7389 | 0.7252 | 0.51 | 0.44 | 0.44 | 0.70 | 0.70 |
| | K-Nearest Neighbors | Aptos | 0.6369 | 0.6987 | 0.43 | 0.44 | 0.42 | 0.65 | 0.67 |
| NEURO-GUARD (Unsupervised Knowledge Extraction) | Logistic Regression | Eye Pacs | 0.6932 | 0.6991 | 0.45 | 0.40 | 0.42 | 0.68 | 0.69 |
| | Random Forest | Eye Pacs | 0.7198 | 0.7342 | 0.50 | 0.44 | 0.47 | 0.73 | 0.74 |
| | SVM | Eye Pacs | 0.7011 | 0.7103 | 0.48 | 0.42 | 0.45 | 0.72 | 0.72 |
| | Gradient Boosting | Eye Pacs | 0.7429 | **0.7465** | 0.55 | 0.48 | 0.52 | 0.75 | 0.75 |
| | K-Nearest Neighbors | Eye Pacs | 0.7124 | 0.7206 | 0.50 | 0.43 | 0.46 | 0.72 | 0.73 |
| NEURO-GUARD (Supervised Artifact Detection) | Logistic Regression | Aptos | 0.7322 | 0.7732 | 0.59 | 0.49 | 0.49 | 0.77 | 0.75 |
| | Random Forest | Aptos | 0.8005 | 0.7978 | 0.65 | 0.56 | 0.57 | 0.80 | 0.78 |
| | SVM | Aptos | 0.7432 | 0.7814 | 0.59 | 0.50 | 0.50 | 0.77 | 0.75 |
| | Gradient Boosting | Aptos | 0.8415 | **0.8469** | 0.69 | 0.58 | 0.58 | 0.83 | 0.80 |
| | K-Nearest Neighbors | Aptos | 0.7896 | 0.7814 | 0.63 | 0.56 | 0.55 | 0.78 | 0.76 |
| NEURO-GUARD (Supervised Artifact Detection) | Logistic Regression | Eye Pacs | 0.7354 | 0.7401 | 0.52 | 0.46 | 0.48 | 0.75 | 0.75 |
| | Random Forest | Eye Pacs | 0.7651 | 0.7713 | 0.58 | 0.51 | 0.54 | 0.78 | 0.78 |
| | SVM | Eye Pacs | 0.7423 | 0.7512 | 0.55 | 0.48 | 0.51 | 0.76 | 0.76 |
| | Gradient Boosting | Eye Pacs | 0.7742 | **0.7796** | 0.60 | 0.53 | 0.56 | 0.80 | 0.80 |
| | K-Nearest Neighbors | Eye Pacs | 0.7545 | 0.7608 | 0.56 | 0.49 | 0.52 | 0.77 | 0.77 |

Table 2. Performance comparison of our framework with DeepX-SOZ, and CNN for seizure onset zone (SOZ) identification. The benchmark results are based on anatomical MRI-based manual SOZ identification and surgical outcomes. Our method achieves **83.27%**, demonstrating strong generalization across multi-center datasets.

| Method | SOZ Identification AUC |
|---|---|
| DeepXSOZ | 81.6 |
| CNN | 46.1 |
| **Our Framework** | **83.27** |

Table 3. Summary of VLM Observations for DR Classification

| Model | Key Observations |
|---|---|
| Qwen-VL-Chat | Generated strong textual reasoning yet achieved only 10.4% accuracy, notably below its 33.84% CARES benchmark. |
| LLaVa-Med-v1.5-mistral-7b | Improved classification accuracy to 24.83% by explicitly defining features, but still prone to hallucination. |
| LLaVa-v1.6-vicuna-7b | Performed admirably in captioning but failed in fine-grained feature detection, limiting reliability for DR classification. |
| Molmo-7B-D-0924 | Demonstrated robust reasoning yet frequently hallucinated features, resulting in misclassifications. |
| CogVLM | Excelled at descriptive captioning capabilities but lacked structured classification functionality. |

Claude, and GPT-4 for structured feature extraction. While LLMs excelled at reasoning, they struggled with direct visual grounding. For instance, Gemini often misclassified bright artifacts as exudates, and MedGPT produced numerous false positives despite domain-specific tuning. GPT-4 achieved the highest accuracy (73%) through iterative bounding-box refinement but remained inconsistent [19].

**Phase 3: The NEURO-GUARD Solution.** The failure of black-box models to provide reliable, non-hallucinated explanations prompted us to move beyond direct prediction. We propose **NEURO-GUARD**, an agentic framework where the LLM does not diagnose directly but instead acts as a logic engine. It generates deterministic Python code (e.g., OpenCV functions) to extract features based on clinical expert rules, ensuring that the diagnosis is grounded in verifiable image data rather than stochastic model outputs.We compare the performance of traditional deep learning models, VLMs, LLMs, and the specialized fine-tuned CLIP-DR model on diabetic retinopathy with the results of our proposed framework, which clearly demonstrates superior outcomes 4.

## 4.1. Localization and Clinical Interpretability Assessment for Diabetic Retinopathy (DR)

To evaluate lesion-level interpretability, we use a YOLOv11 detector fine-tuned on 300 expert-annotated retinal images containing exudates, hemorrhages, cotton-wool spots, and venous abnormalities. The clinical categories used for annotation are first retrieved through our RAG module from PubMed and then verified by ophthalmology experts.

NEURO-GUARD generates the complete YOLOv11 training and inference code through its LLM-driven code synthesis pipeline. The resulting model achieves an 88% localization accuracy, with Intersection-over-Union (IoU) used for automatic verification of predicted lesion regions. This integration of symbolic knowledge, expert annotations, and LLM-generated code enables clinically aligned, pixel-level localization essential for trustworthy diagnosis.

## 5. Results and Comparison

We conducted a comprehensive evaluation of existing Vision-Language Models (VLMs), fine-tuned VLMs, and specialized models against our proposed NEURO-GUARD framework using the same dataset, data processing pipeline, and accuracy metrics to ensure a fair comparison. The baseline reference, a Vision Transformer (ViT)-based model optimized through simultaneous parameter optimization and a feature-weighted ECOC ensemble, achieved an accuracy of 78.4% as reported in IEEE paper. Standard VLMs such as Molmo-7B-D-0924 and LLaVa-v1.6-vicuna-7b demonstrated significantly lower accuracies of 17.22% and 8%, respectively. Medical Fine-Tuned VLMs, including LLaVa-Med-v1.5-mistral-7b and Med Flamingo-9B, showed moderate improvements with accuracies of 24.83% and 5%, respectively. In contrast, the specialized model CLIP-DR achieved a considerably higher accuracy of 69.70%.

Our NEURO-GUARD framework, which integrates domain-specific knowledge and iterative verification mechanisms, outperformed all baseline models. Without incorporating ground truth data of knowledge while learning the knowledge components through open CV and other zeo shot models, NEURO-GUARD achieved an accuracy of 73.84%, closely approaching the specialized CLIP-DR model. When ground truth data of human annotated knowledge components on image was integrated, NEURO-GUARD's accuracy surged to 84.69%, surpassing the ViT benchmark and demonstrating the efficacy of our holistic framework. This significant improvement underscores the advantage of leveraging additional domain-specific steps such as knowledge retrieval and code refinement, which are absent in the baseline VLMs and fine-tuned models. Cross-dataset generalization (train: APTOS → test: EyePACS) improves by 5.2% compared to the ViT baseline, illustrating the value of clinically grounded reasoning for domain-invariant performance.

## 6. Discussion

The experimental results validate NEURO-GUARD's ability to reconcile diagnostic accuracy with clinical interpretability, addressing a core limitation of modern medical AI systems. As shown in Table 1, NEURO-GUARD achieves 84.69% accuracy on the Aptos dataset with super-

vised artifact detection through yolo object detction pipeline to extract knowledge, surpassing the ViT benchmark by 6.2% as shown in Table 4 . This performance gain underscores the value of grounding vision models in domain-specific knowledge a critical factor missing in standard VLMs like LLaVa-v1.6-vicuna-7b (8.6% accuracy) and even specialized models like CLIP-DR (69.7%). The framework's superiority is further evident in its generalization to Eye Pacs (77.96% accuracy) and MRI-based seizure detection (83.27% SOZ accuracy, as shown in Table 2 ), demonstrating robustness across modalities and clinical tasks.

### 6.1. Three key insights emerge:

**Entropic Reward Drives Precision:** NEURO-GUARD's self-verification mechanism (Eq. 3–4) reduces uncertainty in feature extraction, as seen in the 7.85% accuracy jump when integrating ground truth knowledge from expert human annotation of knowledge components (73.84% → 84.69%). Gradient Boosting consistently outperformed other classifiers (Table 1), suggesting that ensemble methods better capture the probabilistic dependencies between LLM-generated features and diagnostic labels.

**VLMs Fail in Fine-Grained Medical Reasoning:** Standard and medical-tuned VLMs (e.g., Med Flamingo-9B at 5.3%) perform poorly due to their reliance on coarse image-text correlations rather than biomarker-level grounding. NEURO-GUARD circumvents this by decomposing diagnoses into executable code (e.g., YOLO lesion detectors), ensuring features align with ICD-10 criteria. However, a known failure mode is LLM hallucinations in code generation, where incorrect feature-detection logic may emerge. To mitigate this, NEURO-GUARD integrates reinforcement learning-based verification, but further improvements such as hybrid clinician-AI oversight may be required in critical applications.

**Generalizability vs. Specialization Trade-Off:** While DeepXSOZ achieves higher SOZ accuracy (81.6%), NEURO-GUARD's 83.27% accuracy with multi-center MRI data highlights its adaptability and contribution towards more improved results. Unlike task-specific models (e.g., CLIP-DR), NEURO-GUARD's modular architecture allows seamless integration of new knowledge graphs, enabling rapid adaptation to evolving clinical guidelines. However, real-world deployments may face challenges in low-resource settings where access to clinician-validated ground truth is limited. Future iterations could incorporate self-supervised learning strategies to improve performance in such scenarios.

A limitation is the dependency on clinician-validated ground truth (Supervised artifact detection required human annotated data) for optimal performance. However, the framework's reinforcement learning pipeline mitigates this by iteratively refining feature extractors using entropic re-

Table 4. Comparison of Accuracy Across Different Models

| Method | Accuracy |
|---|---|
| **ViT Benchmark** | 78.40% |
| **Standard VLMs** | |
| Molmo-7B-D-0924 | 17.22% |
| LLaVa-v1.6-vicuna-7b | 8.60% |
| Qwen-VL-Chat-7b | 10.00% |
| **Medical Fine-Tuned VLMs** | |
| LLaVa-Med-v1.5-mistral-7b | 24.83% |
| Med Flamingo-9B | 5.30% |
| **Specialized Models** | |
| CLIP-DR | 69.70% |
| **Our Framework** | |
| NEURO-GUARD - (Unsupervised) | 73.84% |
| NEURO-GUARD - (Supervised) | **84.69%** |

wards, reducing manual annotation demands over time. Additionally, by aligning extracted features with clinically validated criteria, NEURO-GUARD has the potential to reduce misdiagnosis rates in real-world clinical trials, providing interpretable insights that improve diagnostic trustworthiness and patient outcomes.

## 7. Ethical Considerations and Deployment Challenges

While NEURO-GUARD demonstrates strong interpretability and diagnostic accuracy, real-world deployment in clinical settings requires ethical scrutiny. First, reliance on LLM-generated code introduces risks of hallucinated logic, which could compromise patient safety. Reinforcement learning mitigates some of these risks, but human-in-the-loop validation remains critical, especially in high-stakes settings such as ophthalmology or neurology.

Secondly, access to clinician-validated ground truth is limited in many regions. This may hinder fairness and lead to performance gaps across populations. NEURO-GUARD's self-verification partially addresses this, but further safeguards such as clinical audits, bias analysis, and explainable output formats must be instituted prior to deployment.

Finally, the use of external clinical data (e.g., PubMed) raises data privacy and provenance concerns. Deployment should ensure that all retrieved sources are compliant with medical data standards and not inadvertently expose sensitive patient information

## 8. Conclusion

NEURO-GUARD redefines the paradigm of medical AI by unifying the complementary strengths of self-supervised vision models and LLMs. By formalizing feature extraction as a knowledge-guided code generation task, the framework achieves state-of-the-art accuracy while producing explanations grounded in clinical guidelines. Key innovations include:

Clinically Aligned Interpretability: NEURO-GUARD generates reports that map detected biomarkers (e.g., hemorrhages) to diagnostic criteria, bridging the gap between AI outputs and clinician workflows.

Self-Verification via Entropic Rewards: The LLM-driven reinforcement learning loop ensures feature extractors evolve with medical knowledge, minimizing hallucinations and distributional shifts.

Cross-Modal Generalizability: Validated on diabetic retinopathy (Aptos, Eye Pacs) and seizure detection (MRI), NEURO-GUARD demonstrates versatility across imaging modalities and diseases.

The framework's codebase and models are open-sourced to accelerate research in interpretable medical AI. Future work will extend NEURO-GUARD to video-based diagnostics (e.g., echocardiography) and federated learning for privacy-sensitive deployments.

## 9. Limitations and Future Work

While NEURO-GUARD demonstrates strong interpretability and cross-domain performance, several limitations remain. First, the framework relies on RAG-based retrieval, which can propagate outdated or incomplete clinical knowledge into the rule base, influencing code generation quality. Second, the multi-stage prompting and verification loop introduces computational overhead, limiting real-time deployment. Third, although entropic self-verification reduces hallucinations, LLM-generated code may still exhibit logical inconsistencies in rare cases. Finally, the reliance on clinician-validated ground truth for optimal refinement can restrict performance in low-resource settings.

Future work will address these limitations through a unified medical foundation model that integrates structured clinical ontologies, symbolic reasoning, and pixel-level supervision within a single architecture. Such a model would eliminate external retrieval noise, provide stable lesion-aware representations, and enable end-to-end self-consistency checks without handcrafted prompts. We also plan to incorporate continual learning to update clinical knowledge dynamically, and explore on-device optimization for real-time diagnostic support in resource-constrained environments by aiming to minimize hallucination of Large foundational models.

# References

[1] F. Antaki et al. Vision-language models for feature detection of macular diseases on optical coherence tomography. *JAMA Ophthalmology*, 142(4), 2024. 5

[2] B. Cohen-Wang et al. Trustworthy medical imaging with large language models: A study of hallucinations. *arXiv preprint arXiv:2408.xxxxx*, 2024. 4

[3] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, others, and N. Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *Proceedings of ICLR*, 2021. 1, 3

[4] Z. Gu, C. Yin, F. Liu, and P. Zhang. Medvh: Towards systematic evaluation of hallucination for large vision language models in the medical context. *arXiv preprint arXiv:2407.02730*, 2024. 5

[5] T. Gupta and A. Kembhavi. Visual programming: Compositional visual reasoning without training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 4

[6] R. Holland et al. Medical hallucination in foundation models and their impact on healthcare. *arXiv preprint arXiv:2402.XXXXX*, 2024. 5

[7] P. W. Koh, T. Nguyen, Y. S. Tang, S. Mussmann, E. Pierson, B. Kim, and P. Liang. Concept bottleneck models. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2020. 3

[8] S. M. Lundberg and S.-I. Lee. A unified approach to interpreting model predictions. In *Proceedings of NeurIPS*, 2017. 1

[9] P. Mlynarski, H. Chen, and L. Huang. Object detection-based approach for medical image classification and reasoning. *Medical Image Analysis*, 92:102658, 2023. 5

[10] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, others, and D. Hassabis. Human-level control through deep reinforcement learning. *Nature*, 2015. 2

[11] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever. Learning transferable visual models from natural language supervision. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2021. 3

[12] C. Rudin. Stop explaining black-box machine learning models for high-stakes decisions and use interpretable models instead. *Nature Machine Intelligence*, 1(5):206–215, 2019. 3

[13] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of ICCV*, 2017. 1, 3

[14] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *Proceedings of ICLR*, 2015. 1, 3

[15] A. Singhal et al. Advancements in large language models for medical diagnostics. In *Proceedings of EMNLP*, 2023. 3

[16] D. Surís, S. Menon, and C. Vondrick. Vipergpt: Visual inference via python execution for reasoning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023. 4

[17] Y. Tang, X. Li, and M. Zhang. Few-shot object detection for lesion localization in retinal images. *IEEE Access*, 10: 91245–91257, 2022. 5

[18] M. Volpi, Z. Zhang, and Y. Chen. Robustness of deep learning models in medical imaging: A survey. In *Proceedings of MICCAI*, 2018. 2

[19] X. Wang et al. Llm-seg: Bridging image segmentation and large language model reasoning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 6

[20] Y. Wang, K. Xu, and G. Lu. Enhancing medical image classification with domain-specific knowledge integration. In *Proceedings of ICML*, 2021. 2

[21] Z. Wang, Z. Wu, D. Agarwal, and J. Sun. Medclip: Contrastive learning from unpaired medical images and text. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2022. 3

[22] Y. Wu, T. Zhang, and J. Holmes. Reinforcement learning for interpretable medical image classification. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, 2022. 1, 3

[23] P. Xia, Z. Chen, J. Tian, Y. Gong, et al. Cares: A comprehensive benchmark of trustworthiness in medical vision language models. In *Proceedings of the Neural Information Processing Systems (NeurIPS)*, 2024. 5

[24] A. Yan, Y. Wang, Y. Zhong, Z. He, P. Karypis, Z. Wang, C. Dong, A. Gentili, C. N. Hsu, J. Shang, and J. J. McAuley. Robust and interpretable medical image classifiers via concept bottleneck models. *arXiv preprint arXiv:2310.03182*, 2023. 3

[25] Y. Zhang et al. Samed: Adapting segmentation algorithms to medical imaging. In *Proceedings of MICCAI*, 2022. 4

[26] X. Zhou, Y. Li, and P. Chen. Towards interpretable medical imaging with deep learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 1