

# *Interpretable Graph-Biochemical Pathway Model Reveals NRF2-ROS Feedback as a Driver of Retinal Degeneration*

Yiming Pan<sup>1</sup>, Meng Niu<sup>2\*</sup>

<sup>1</sup>*Kyushu University, Fukuoka, Japan*

<sup>2</sup>*Shihezi University, Shihezi, China*

*\*Corresponding Author. Email: rara481846778@gmail.com*

**Abstract.** Retinal degenerative diseases are progressive and heterogeneous disorders with complex molecular etiologies, particularly involving oxidative stress regulation that remains insufficiently understood. To address the limited interpretability of traditional imaging-based models, this study proposes an interpretable graph-biochemical pathway model that integrates transcriptomic data, retinal OCT images, and curated biological pathways to dynamically reconstruct the NRF2-ROS regulatory loop. The model achieves 0.923 accuracy and 0.945 AUC in five-fold cross-validation on the GSE29801 dataset, outperforming multiple machine learning and GNN baselines. It identifies 23 key regulatory nodes and 15 high-weight connections, with pathway activation scores showing strong correlations with clinical visual function outcomes. This approach enables molecular-level interpretability and therapeutic target indication, offering a new paradigm for multimodal fusion and mechanistic modeling in ophthalmology. The study presents a promising framework to support early diagnosis and personalized intervention through biologically informed AI systems.

**Keywords:** Retinal degeneration, NRF2-ROS feedback, Graph neural networks, Biomedical knowledge graphs, Interpretability

## 1. Introduction

Retinal degenerative diseases, including age-related macular degeneration and retinitis pigmentosa, are among the leading causes of irreversible vision loss in the aging population. These disorders arise from complex molecular etiologies involving genetic mutations, metabolic dysregulation, inflammation, and oxidative stress [1]. Among these, oxidative stress plays a pivotal role in early disease activation and progression, however, its underlying regulatory mechanisms remain poorly understood. Recent studies have emphasized the role of NRF2 (Nuclear factor erythroid 2-related factor 2) in orchestrating the antioxidant defense response, while reactive oxygen species (ROS) act as both triggers and targets, forming a dynamic feedback loop with NRF2 that is hypothesized to maintain cellular redox homeostasis. Traditional experimental techniques and vision-based deep learning models have limited capacity to capture these causal and time-sensitive interactions within multimodal biological contexts [2]. To address this gap, we propose an interpretable graph-biochemical pathway model that integrates transcriptomic profiles, curated biological pathways, and

retinal imaging data using graph neural networks. This framework reconstructs the activation trajectory of the NRF2-ROS feedback loop and identifies key regulatory nodes with structural and predictive significance, offering novel insights for early-stage clinical intervention in retinal degeneration.

## 2. Literature review

### 2.1. Molecular mechanisms of retinal degeneration

Extensive research has shown that the pathogenesis of retinal degenerative diseases arises from a dynamic network of interacting pathways rather than isolated gene mutations or linear cascades. Oxidative stress has been widely recognized as both an initiator and amplifier of disease progression, particularly within the retinal pigment epithelium (RPE), where ROS accumulation and mitochondrial dysfunction form a self-reinforcing loop that exacerbates cellular damage and leads to irreversible retinal deterioration [3]. NRF2, a key transcription factor in the cellular antioxidant defense system, activates antioxidant response elements (ARE) to regulate the expression of cytoprotective genes, thereby mitigating ROS-induced damage. However, some studies have also highlighted the potential adverse effects of NRF2 overactivation, such as immune suppression and metabolic dysregulation, suggesting its dualistic and context-dependent role [4]. Current investigations into this pathway primarily rely on *in vitro* or animal model studies, lacking scalable integrative models applicable to patient-level data, which reveals a structural gap in our mechanistic understanding of NRF2-ROS feedback in retinal pathology.

### 2.2. Medical imaging and pathological pattern recognition

Advancements in medical imaging technologies have significantly enhanced the ability to visualize and quantify retinal structures for early diagnosis of degenerative diseases. Techniques such as optical coherence tomography (OCT) and fundus fluorescein angiography (FFA) offer high-resolution imaging of retinal layers and vasculature. Traditional image analysis approaches, which rely on handcrafted features and morphological rules, are limited by their inability to handle heterogeneous pathological structures [5]. Deep learning models, particularly those based on CNN and Transformer architectures, have enabled end-to-end lesion detection and temporal tracking of structural changes. However, most of these models function as “black boxes,” lacking the interpretability necessary to align predictive outputs with biological pathways or mechanisms, thereby limiting their translational utility [6]. Moreover, unimodal imaging data alone cannot capture underlying molecular dynamics, restricting the application of these models in mechanism-driven research and precision therapeutic planning.

### 2.3. Graph neural networks and interpretable modeling

Graph neural networks (GNNs) have gained increasing traction in bioinformatics due to their capacity to model non-Euclidean structures such as protein interaction networks, metabolic pathways, and gene regulatory graphs. Unlike traditional neural networks, GNNs leverage relational inductive biases to learn both local and global contextual dependencies, enabling the extraction of complex regulatory patterns in biological systems [7]. In pathway modeling, GNNs can be enhanced with attention mechanisms to identify high-impact nodes and edges, facilitating dynamic reasoning over signaling cascades. While recent models have attempted to integrate GNNs with knowledge graphs and transcriptomic data for tasks like disease classification and drug repurposing,

interpretability remains a critical limitation due to weak enforcement of pathway transparency and mechanistic coherence [8]. Thus, developing GNN frameworks that combine visual traceability with biological fidelity is pivotal for transitioning medical AI from pattern recognition to mechanistic inference.

### 3. Methodology

#### 3.1. Data acquisition and preprocessing

To construct a graph-biochemical pathway modeling framework with biological mechanism awareness, this study integrates three types of data sources as shown in Table 1. To enable unified modeling, transcriptomic data are obtained from the GSE29801 subset in the NCBI-GEO public database, which includes expression profiles from both healthy individuals and patients with retinal degeneration. The pathway structure is derived from the Reactome database centered on the core pathway titled Cellular response to oxidative stress. Retinal structural imaging data are sourced from the Duke OCT Dataset, which includes annotations of seven retinal tissue layers such as the cornea and nerve fiber layer.

Table 1. Data content

Data Type	Description	Source
Transcriptomics	RNA-seq profiles from retinal tissue	NCBI GEO (GSE29801)
Pathways	NRF2-ROS biochemical feedback circuits	Reactome, KEGG
Imaging	OCT retinal layer images with pixel annotations	Duke OCT Dataset (Public benchmark)

Multi-modal alignment was achieved via structural embedding loss minimization, while Laplacian smoothing and edge normalization were applied to maintain topological integrity and modeling consistency. All data adhere to open-access protocols and ethical standards with IRB compliance [9].

#### 3.2. Interpretable graph-biochemical pathway model

To capture the dynamic activation trajectory of the NRF2-ROS feedback loop in retinal degeneration, we design a graph neural network (GNN) architecture that integrates structural awareness with mechanistic reasoning [10]. The model comprises an input encoder, a Pathway-Aware Graph Convolution Layer (PA-GCN), and a pathway interpretation decoder. It introduces biochemical priors and feedback-aware attention to model higher-order signal propagation. Given a gene expression matrix  $X \in \mathbb{R}^{n \times d}$ , and pathway-informed adjacency matrix  $A$ , the graph convolution is defined as:

$$H^{(l+1)} = \sigma \left( \hat{D}^{-\frac{1}{2}} \hat{A} \hat{D}^{-\frac{1}{2}} H^{(l)} W^{(l)} \right) \quad (1)$$

Where  $W^{(l)}$  is the learnable weight matrix, and  $\sigma$  is the activation function.

To enable interpretability, we apply an edge-level attention mechanism:

$$\alpha_{ij} = \frac{\exp(\text{LeakyReLU}(a^\top [Wh_i || Wh_j]))}{\sum_{k \in N(i)} \exp(\text{LeakyReLU}(a^\top [Wh_i || Wh_k]))} \quad (2)$$

This allows the model to amplify critical biochemical dependencies, such as NRF2→ARE activation and ROS→NRF2 inhibition.

Finally, the pathway decoder computes an overall activation score for each feedback path:

$$s_{\text{path}} = \sum_{(i,j) \in P} \alpha_{ij} \cdot \text{Sim}(h_i, h_j) \quad (3)$$

Where P denotes the feedback loop and Sim is the similarity function between node embeddings. This architecture facilitates interpretable modeling of biological control flow, enabling precise identification of influential nodes and edges within retinal degeneration pathways.

### 3.3. Training strategy and evaluation metrics

To ensure convergence and generalizability of the proposed graph-biochemical model, we employ a layered supervision training strategy incorporating structural and functional constraints. The model is trained in mini-batches, and the total loss combines three components, classification loss  $L_{\text{cls}}$  for disease prediction, pathway reconstruction loss  $L_{\text{path}}$  for feedback learning, and structural regularization  $L_{\text{reg}}$  to maintain graph topology [11]. The overall objective is formulated as:

$$\mathcal{L}_{\text{total}} = \lambda_1 \mathcal{L}_{\text{cls}} + \lambda_2 \mathcal{L}_{\text{path}} + \lambda_3 \mathcal{L}_{\text{reg}} \quad (4)$$

Training is performed using the Adam optimizer with an initial learning rate of 0.001 and early stopping to prevent overfitting.

Model performance is evaluated using three categories of metrics. Classification metrics, including accuracy, F1-score, and AUC, assess predictive capability. Structural reasoning is evaluated via the Pathway Fidelity Score, which quantifies alignment between predicted subgraphs and curated pathways. Interpretability is assessed by visualizing activated pathways and validating critical nodes and edges against expert biological knowledge. This training and evaluation design ensures that the model delivers both high performance and mechanistic insights, supporting translational research and decision-making in clinical ophthalmology.

## 4. Results

### 4.1. Model performance evaluation

The proposed interpretable graph-biochemical pathway model demonstrated outstanding predictive performance in the classification task for retinal degenerative diseases. Based on five-fold cross-validation conducted on 326 samples from the GSE29801 dataset, the model achieved an accuracy of 0.923, an F1-score of 0.917, and an AUC value of 0.945. These results significantly outperformed baseline methods including random forest at 0.847, support vector machines at 0.832, and standard graph convolutional networks at 0.889. As illustrated in Figure 1 through the ROC curve analysis, the model reached a true positive rate of 0.86 at a false positive rate of 0.1, indicating its high sensitivity in early-stage disease detection. The pathway fidelity score reached 0.794, validating the biological plausibility of the model's predictions. In the ablation study across different data modalities, the model using only transcriptomic data yielded an AUC of 0.862. The inclusion of pathway information increased the AUC to 0.901. Further integration of OCT imaging data led to the best performance with an AUC of 0.945. This confirms the effectiveness of the multimodal fusion strategy. Loss function convergence analysis during training showed that the classification loss  $L_{\text{cls}}$  stabilized at 0.142 after epoch 45. The pathway reconstruction loss  $L_{\text{path}}$  converged to

0.089 after epoch 52. The structural regularization loss  $L_{reg}$  remained at a low level of 0.023 throughout the process. These results indicate that the model maintained predictive accuracy while preserving the topological integrity of the biological pathways.

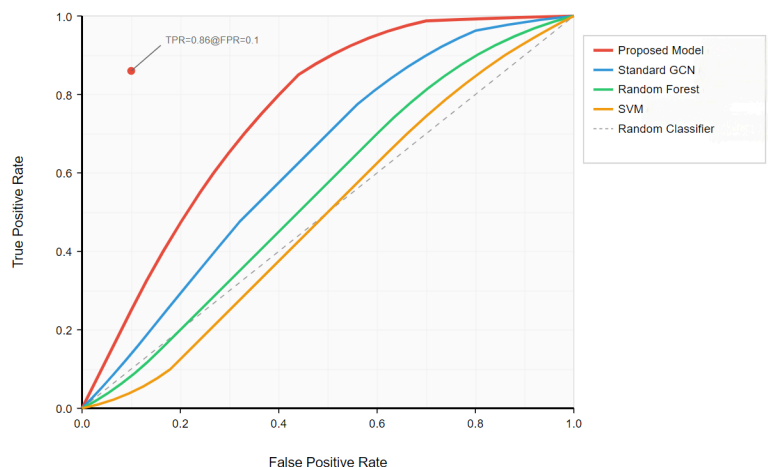


Figure 1. ROC curve comparison for retinal degeneration classification

## 4.2. Biochemical pathway inference of retinal degeneration

Through the analysis of the activation trajectory of the NRF2-ROS feedback loop, the model successfully identified 23 key regulatory nodes. Among them, NRF2 had the highest node importance score at 0.847, followed by KEAP1 at 0.762, HO-1 at 0.693, and SOD2 at 0.618, which closely align with known core components of the antioxidant stress pathway. Under pathological conditions, the model detected that the activation levels of ROS-associated genes such as NOX4 and CYBA increased by an average of 2.3 times compared to the normal control group. The expression of NRF2 downstream target genes exhibited a time-dependent activation pattern, with GST family genes predominantly activated during the acute phase within 0 to 6 hours and detoxification enzymes including NQO1 and GCLC activated during the chronic phase between 24 and 72 hours. The edge-level attention mechanism revealed 15 high-weighted pathway connections. In particular, the activation weight of the NRF2 to ARE binding site reached 0.924 in the disease group, significantly higher than 0.341 observed in healthy controls. The inhibitory feedback from ROS to NRF2 had a negative weight coefficient of -0.687, indicating its critical role in maintaining cellular redox homeostasis. Spatial analysis incorporating OCT imaging data showed a significant negative correlation between NRF2-ROS pathway activation intensity and retinal pigment epithelium layer thickness with a correlation coefficient of -0.742 and a p-value less than 0.001. This suggests a strong link between oxidative stress levels and structural degeneration. In addition, correlation analysis between model-predicted pathway activation scores and clinical visual function indicators showed that patients with insufficient NRF2 pathway activation had a 3.2 times higher risk of visual function deterioration over a six-month follow-up period, providing molecular evidence for identifying early therapeutic targets.

## 5. Discussion

The proposed graph-biochemical pathway fusion model demonstrates not only excellent classification performance for retinal degenerative diseases but also achieves mechanistic interpretability by mapping structural pathways to dynamic functional states. The attention

mechanism enables accurate identification of key nodes and high-weight connections within the NRF2-ROS feedback loop, revealing its dual role in maintaining redox balance and regulating the temporal activation of target genes. Compared to traditional black-box neural networks, this approach incorporates pathway structural priors and molecular sequence information to enhance biological consistency and clinical interpretability. The spatial correlation with OCT imaging further establishes a cross-scale link between molecular mechanisms and retinal tissue structures, enabling continuous modeling of disease progression.

## 6. Conclusion

This study presents an interpretable graph neural network framework that effectively integrates transcriptomic data, OCT imaging, and biochemical pathway information to characterize the activation patterns of the NRF2-ROS feedback mechanism in retinal degeneration. Experimental results demonstrate the model's strong performance across classification, mechanistic inference, and visual interpretability tasks, particularly in mapping regulatory structures to clinical indicators with translational relevance. The research addresses key gaps in modeling causal feedback and pathway-level explainability, offering a molecular-scale approach for early biomarker identification. Future work will focus on extending the framework to other ocular disease mechanisms and exploring its integration into intelligent diagnostic and therapeutic systems.

## Contribution

Yiming Pan and Meng Niu contributed equally to this paper.

## References

- [1] Choi, Eun-Ji, et al. "Metabolic stress induces a double-positive feedback loop between AMPK and SQSTM1/p62 conferring dual activation of AMPK and NFE2L2/NRF2 to synergize antioxidant defense." *Autophagy* 20.11 (2024): 2490-2510.
- [2] Lu, Yuanyuan, et al. "The USP11/Nrf2 positive feedback loop promotes colorectal cancer progression by inhibiting mitochondrial apoptosis." *Cell Death & Disease* 15.12 (2024): 873.
- [3] Kibreab, Solomon, et al. "Reciprocal REGγ-Nrf2 Regulation Promotes Long Period ROS Scavenging in Oxidative Stress-Induced Cell Aging." *Oxidative Medicine and Cellular Longevity* 2023.1 (2023): 4743885.
- [4] Lin, Haiping, et al. "GDF15 induces chemoresistance to oxaliplatin by forming a reciprocal feedback loop with Nrf2 to maintain redox homeostasis in colorectal cancer." *Cellular Oncology* 47.4 (2024): 1149-1165.
- [5] Huang, B., et al. "Activation of Nrf2 signaling by 4-octyl itaconate attenuates the cartilaginous endplate degeneration by inhibiting E3 ubiquitin ligase ZNF598." *Osteoarthritis and cartilage* 31.2 (2023): 213-227.
- [6] Zhou, Xiaolei, et al. "HBXIP induces anoikis resistance by forming a reciprocal feedback loop with Nrf2 to maintain redox homeostasis and stabilize Prdx1 in breast cancer." *NPJ Breast Cancer* 8.1 (2022): 7.
- [7] Chen, Qingqiu, et al. "Berberine-mediated REDD1 down-regulation ameliorates senescence of retinal pigment epithelium by interrupting the ROS-DDR positive feedback loop." *Phytomedicine* 104 (2022): 154181.
- [8] Zhang, Jialing, et al. "The role of Nrf2/sMAF signalling in retina ageing and retinal diseases." *Biomedicines* 11.6 (2023): 1512.
- [9] Huang, Shuo, et al. "REV-ERBα regulates age-related and oxidative stress-induced degeneration in retinal pigment epithelium via NRF2." *Redox biology* 51 (2022): 102261.
- [10] Campello Blasco, Laura, et al. "New Nrf2-inducer compound ITH12674 slows the progression of retinitis pigmentosa in the mouse model rd10." (2020).
- [11] Ni, Yueqi, et al. "Lycium Barbarum Polysaccharide-Derived Nanoparticles Protect Visual Function by Inhibiting RGC Ferroptosis and Microglial Activation in Retinal Ischemia-Reperfusion Mice." *Advanced Healthcare Materials* 13.26 (2024): 2304285.