

Diabetic Retinopathy: A Proposed Multi-Task Deep Learning Model for Lesion Segmentation, Grading, and Localization

Dhruv Premani
B.Tech CSE
Ahmedabad University
Ahmedabad, India
AU2240239

Hariohm Bhatt
B.Tech CSE
Ahmedabad University
Ahmedabad, India
AU2240085

Raj Koticha
B.Tech CSE
Ahmedabad University
Ahmedabad, India
AU2240024

Abstract—Diabetic retinopathy (DR) is a complex disease that requires precise lesion segmentation, accurate disease grading, and reliable anatomical localization. In this report, we propose a novel multi-task deep learning model that integrates convolutional neural networks (CNNs), Transformer-based attention, and multi-task learning to jointly address segmentation, grading, and localization in fundus images. Our framework leverages advanced pretraining, data augmentation, and task-specific fusion modules to capture both local details and global context. We expect that this approach will achieve improvements in segmentation accuracy, grading performance, and localization precision compared to conventional single-task methods.

Index Terms—Diabetic Retinopathy, Multi-Task Learning, Segmentation, Classification, Localization, CNN, Transformer, Deep Learning.

I. INTRODUCTION

Diabetic retinopathy is a leading cause of blindness and requires early, accurate diagnosis to prevent irreversible vision loss. Traditional methods often treat lesion segmentation, disease grading, and anatomical localization as separate tasks. However, such separation neglects the interdependencies among these tasks. Recent studies have shown that integrating these tasks can enhance diagnostic performance by leveraging shared features [1]– [3].

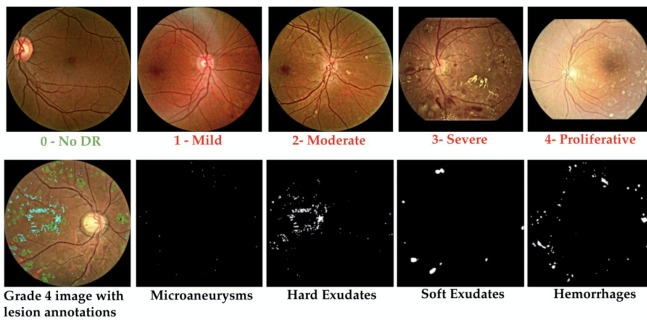


Fig. 1. Example images showing different DR grades (top row: 0–4) and annotated lesions (bottom row). From left to right: 0 (No DR), 1 (Mild), 2 (Moderate), 3 (Severe), and 4 (Proliferative). Bottom images illustrate microaneurysms, hard exudates, soft exudates, and hemorrhages.

In this report, we introduce our proposed model, a unified deep learning approach that combines segmentation, classification, and localization in a single end-to-end trainable network. By fusing CNN and Transformer-based components, the model captures both high-resolution details and global context, which is essential for recognizing subtle lesions and complex disease patterns.

II. METHODOLOGY

Dataset Description

The experiments are conducted on the Indian Diabetic Retinopathy Image Dataset (IDRiD). IDRiD is a public dataset of retinal fundus photographs aimed at developing automated diabetic retinopathy analysis. It is the first database representative of an Indian population with comprehensive annotations, containing typical diabetic retinopathy lesions and normal retinal structures annotated at the pixel level. The dataset is organized into three parts serving different tasks: (A) Segmentation, (B) Disease Grading, and (C) Localization.

For the segmentation task, IDRiD provides 81 color fundus images (split into 54 training and 27 testing images) with expert-drawn pixel-level masks for lesions and anatomical structures. Specifically, each image is annotated with five classes of targets: Microaneurysms (MA), Haemorrhages (HE), Hard Exudates (EX), Soft Exudates (SE), and the Optic Disc (OD). These represent the small reddish microaneurysms, larger blot hemorrhages, bright lipid exudates, pale cotton-wool spots, and the normal optic disc, respectively. A background class (non-lesion) is implicitly present, making it a multi-class segmentation with 6 classes (5 target classes + background). IDRiD is unique in that it includes annotations for both pathology and normal structures in the retina, enabling robust segmentation evaluation.

In addition to segmentation masks, the dataset provides labels for diabetic retinopathy (DR) severity grading—a classification task with 516 images (413 training, 103 testing) labeled with DR grade (0 through 4) and diabetic macular edema (DME) grade (0 through 2) [4]–[6]. This dual design allows researchers to tackle both segmentation and

classification of diabetic retinopathy from the same images, facilitating multi-task learning and comprehensive evaluation of algorithms. All images are high-resolution color fundus photographs (approximately 4288×2848 pixels, 50° field of view) captured under standardized conditions, ensuring sufficient detail for fine lesion segmentation. The availability of both segmentation masks and grading labels makes IDRiD a valuable benchmark for developing and validating retinal image analysis methods. The dataset was originally released as part of an ISBI 2018 challenge on “Diabetic Retinopathy Segmentation and Grading”.

Data Preprocessing

Prior to training, all images and annotation masks undergo a carefully designed preprocessing pipeline. Image resizing is first applied to ensure a uniform input dimension for the network. The original fundus images, which are very high resolution, are scaled down to a manageable size (e.g., 256×256 pixels in our implementation) while approximately preserving the aspect ratio. This reduces memory usage and speeds up training, yet retains essential retinal structures. We choose a resolution that balances computational efficiency with sufficient detail to detect small lesions like microaneurysms. In addition, we crop or pad images as needed to maintain consistency in dimensions across the dataset.

Next, we perform contrast enhancement using Contrast Limited Adaptive Histogram Equalization (CLAHE) on the luminance channel of the images. Specifically, each RGB fundus image is converted from RGB color space to the Lab color space, and CLAHE is applied to the L (lightness) channel. This technique locally enhances the brightness and contrast of retinal images while preserving the original color information (since the a and b color channels remain unchanged). By equalizing illumination in this way, we mitigate lighting variations across the fundus and highlight subtle lesions against the background. After CLAHE, the image is converted back to RGB. This step has the effect of making

small red lesions (MA/HE) and low-contrast lesions (SE) more visible.

Following CLAHE, we apply intensity normalization by scaling pixel values to a $[0, 1]$ range (and optionally standardizing them using the dataset mean and standard deviation). This normalization ensures that all images have comparable intensity statistics, which helps stabilize network training. Overall, the preprocessing pipeline (resizing, CLAHE-based contrast enhancement, and normalization) standardizes the images and improves the visibility of pathological features, providing the model with inputs of consistent quality.

Mask preprocessing is also crucial for multi-class segmentation. The ground truth provided by IDRiD consists of five separate binary masks for each image—one for each lesion type (MA, HE, EX, SE) and one for the optic disc. We merge these into a single combined mask image per fundus. This is done by assigning a unique label value to each mask and overlaying them. For each pixel, if it belongs to a given lesion’s binary mask, that pixel is encoded with the corresponding class index in the combined mask. For example, we assign label 1 for MA, 2 for HE, 3 for EX, 4 for SE, 5 for OD, and use 0 for background. The result is a multi-class segmentation mask where each pixel’s integer value (0–5) indicates the category of that region.

We ensure that there are no overlapping annotations—in the IDRiD data, lesion regions are disjoint by definition (e.g., a pixel cannot simultaneously be labeled as both hemorrhage and exudate). If any overlap were to occur, a priority order would be defined, but in practice, the provided annotations did not conflict. The masks are then resized to the same resolution as the input images (using nearest-neighbor interpolation to preserve label integrity) so that each input image has a corresponding ground truth mask of the same dimensions. By combining the separate binary masks into one, we facilitate training a single multi-class segmentation model that predicts all lesion types in one go, rather than training separate models

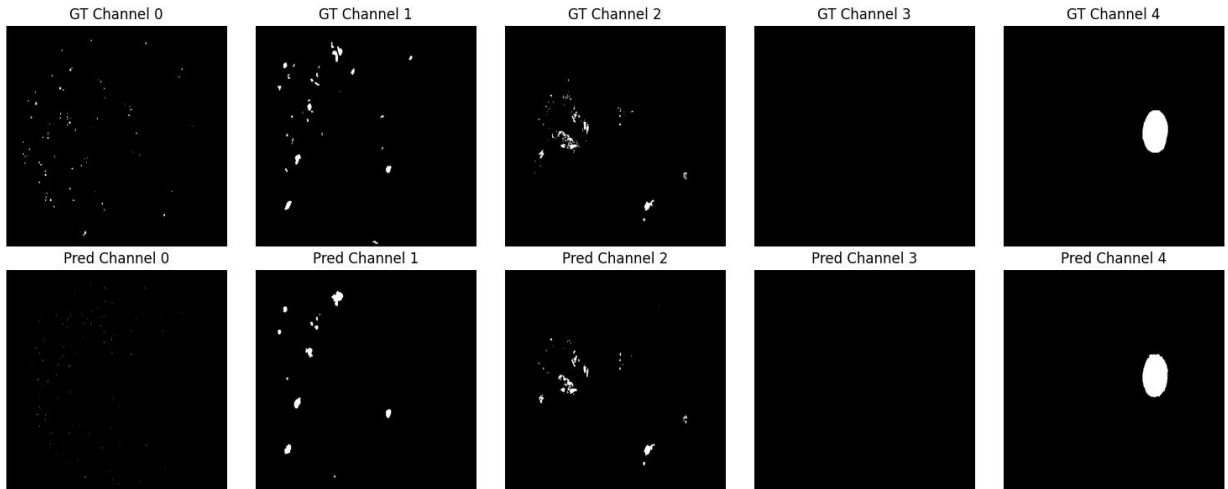


Fig. 2. Segmentation outputs on the testing set. The overlay of predicted (colored) masks with the original fundus images demonstrates the network’s ability to capture both prominent and subtle lesions.

for each lesion.

Model Architecture

We employ an improved U-Net convolutional neural network as the backbone for multi-class segmentation. The architecture is based on the U-Net framework [1], which is an encoder-decoder network with symmetric skip connections. In the encoder path, the image is progressively downsampled through multiple levels: each level consists of convolutional layers that extract feature maps, followed by a downsampling operation (pooling or strided convolution) that reduces spatial dimensions. This yields a hierarchy of feature maps capturing context at increasing scales. At the lowest resolution (the bottleneck), the network has a rich representation of high-level features.

The decoder path then mirrors the encoder: at each level, an upsampling (for example, a transpose convolution) is applied to increase the resolution of the feature maps, and these are concatenated with the corresponding high-resolution feature maps from the encoder via skip connections [1], [2]. These skip connections provide fine-grained spatial detail from earlier layers to the decoder, helping the network accurately localize small structures that might have been lost in downsampling. After concatenation, a series of convolutions produces refined feature maps at that level, and the process repeats until the original image resolution is restored. The final layer is a 1×1 convolution that outputs the desired number of classes (in this case, 6 channels, corresponding to the 5 lesion/OD classes plus background), followed by a softmax activation to produce a probability map for each class.

Our U-Net variant integrates residual learning into the convolutional blocks, inspired by the success of ResNet-like architectures. In a standard U-Net, each stage of the encoder/decoder uses a Double Conv block (two successive 3×3 convolutions with ReLU activations). We replace this with a Residual Double Convolution (ResidualDoubleConv) block. In this design, the input of the block is added to its output, creating a shortcut connection that bypasses the two convolution layers. Specifically, the block performs two convolutions (each followed by a non-linear activation, e.g., ReLU), and then the original input tensor is added element-wise to the final output of these operations.

If the convolutional block changes the number of channels, a 1×1 convolution on the input (the shortcut) is used to match dimensions before addition. This residual addition yields an identity mapping that the conv layers only need to learn a residual of, which significantly eases training of deeper networks. The motivation is that residual connections improve the flow of gradients during backpropagation and help avoid vanishing gradients, allowing us to stack more layers without degradation of performance.

In our architecture, every level’s double-conv module (both in encoder and decoder) is implemented as a ResidualDoubleConv block. Thus, one can view the network as a Residual U-Net (ResUNet)—it preserves the U-shaped encoder-decoder macro-architecture of U-Net but with internal skip connections

within each block. We also include dropout regularization in the model to combat overfitting, given the relatively small training set (54 images). Dropout layers (with a dropout probability of 0.3) are inserted at certain points in the network (such as after the bottleneck or in the decoder) to randomly deactivate a fraction of feature channels during training.

Furthermore, all convolutional layers were initialized with Kaiming He initialization to facilitate faster convergence. Kaiming initialization draws weights from a scaled normal distribution suited for ReLU activations, ensuring that signals propagate without diminishing or exploding in the early stages of training. We did not use batch normalization in this implementation, to avoid issues on small batch sizes; instead, the combination of residual connections and careful initialization was relied upon to stabilize training.

Loss Function

Training the multi-class segmentation model involves optimizing a combined loss function that addresses both class imbalance and region overlap. We use a weighted sum of Focal Loss and Dice Loss as the objective. Focal Loss is an extension of the standard cross-entropy loss designed to handle severe class imbalance by focusing training on hard examples [3]. In fundus lesion segmentation, the imbalance is extreme—lesion pixels (especially microaneurysms) are very sparse compared to background pixels.

Focal Loss introduces a modulating factor $(1 - p_t)^\gamma$ to the cross-entropy, where p_t is the model-predicted probability for the true class, and γ is a focusing parameter [3]. This factor down-weights well-classified pixels and up-weights poorly classified ones. We set $\gamma = 2$ in our experiments, a common choice to strongly focus on harder examples. Additionally, Focal Loss includes a weighting term α to balance classes [3].

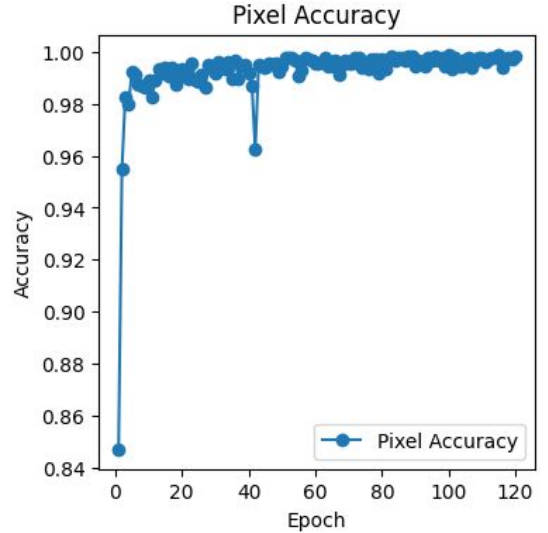


Fig. 3. Pixel accuracy progression during training. The model steadily converges towards high pixel accuracy, demonstrating robust segmentation performance.

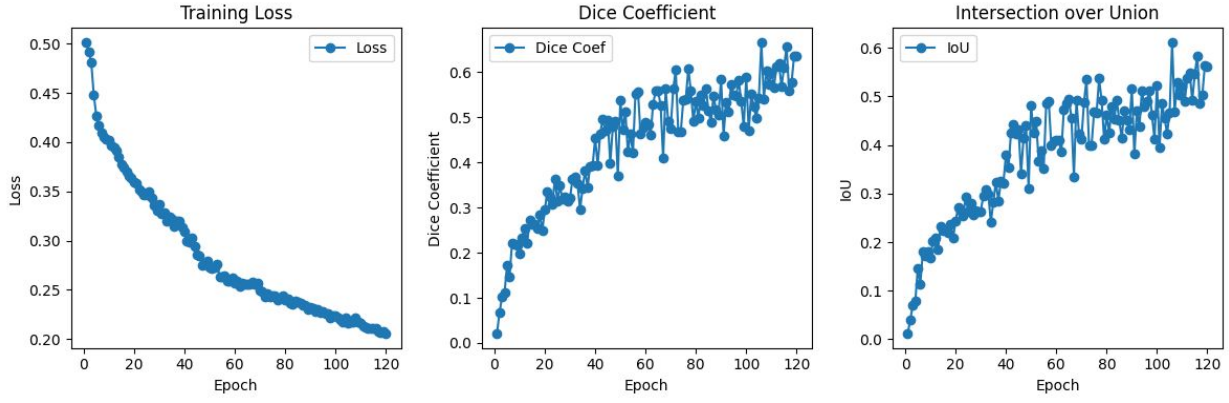


Fig. 4. Evolution of training and validation loss across epochs. The significant reduction in loss reflects effective optimization of our composite Focal + Dice loss function.

We assign a higher weight to lesion classes versus background; for instance, α for lesion pixels might be 0.75 (and 0.25 for background).

While Focal Loss operates at the pixel level, we also include Dice Loss to directly optimize the overlap of predicted masks with the ground truth. Dice Loss is derived from the Dice coefficient, which measures the similarity between two sets. For a given class, the Dice coefficient D between the predicted mask P and ground truth mask G is defined as:

$$D = \frac{2|P \cap G|}{|P| + |G|} = \frac{2TP}{2TP + FP + FN}, \quad (1)$$

where TP , FP , FN are the counts of true positive, false positive, and false negative pixels for that class. Dice Loss is typically defined as $1 - D$, which increases when the prediction and ground truth diverge.

In our multi-class scenario, we compute a Dice Loss for each class and average them (excluding the background class to focus on lesion performance). The total loss \mathcal{L} is:

$$\mathcal{L} = \mathcal{L}_{\text{focal}} + \lambda \mathcal{L}_{\text{dice}}, \quad (2)$$

where we initially set $\lambda = 1$. This composite loss leverages the strengths of both components: the Focal term mitigates class imbalance, while the Dice term ensures good region overlap [3], [6].

Evaluation Metrics

We evaluate segmentation performance using Dice Coefficient, Intersection over Union (IoU), and Pixel Accuracy. For each lesion class c , the Dice Coefficient (DSC) is:

$$\text{DSC}_c = \frac{2TP_c}{2TP_c + FP_c + FN_c}, \quad (3)$$

computed per class and averaged across the five lesion/OD classes. IoU, or Jaccard Index, is:

$$\text{IoU}_c = \frac{TP_c}{TP_c + FP_c + FN_c}, \quad (4)$$

and we report mean IoU (mIoU) over lesion classes. Pixel Accuracy is the fraction of correctly classified pixels:

$$\text{Pixel Accuracy} = \frac{\# \text{ of correctly predicted pixels}}{\# \text{ of total pixels}}. \quad (5)$$

We compute these metrics per batch and epoch, emphasizing Dice and IoU for assessing lesion delineation [2], [3].

Training Protocol

The model is trained using the Adam optimizer with an initial learning rate of 1×10^{-3} and default momentum parameters ($\beta_1 = 0.9$, $\beta_2 = 0.999$). We train for 50 epochs on the 54 training images, with a batch size of 4. Each epoch involves shuffling the training data, computing the combined Focal + Dice loss, and updating weights via backpropagation with gradient clipping (norm threshold 2.0). After each epoch, we evaluate on the test set, saving checkpoints if validation Dice improves. We track training and validation loss, Dice, and IoU, selecting the model with the best validation performance.

Visualization

We visualize segmentation results by overlaying predicted and ground truth masks on test fundus images, using distinct colors for each class (e.g., red for hemorrhages, yellow for exudates). Side-by-side plots compare predicted and ground truth masks, revealing alignment and errors. Visualizations confirm accurate segmentation of large lesions and the optic disc, with minor discrepancies for small lesions like microaneurysms.

Suggestions for Improvement

- **Data Augmentation:** Apply random rotations, flips, scaling, and color jitter to enhance generalization and robustness to imaging variations.
- **Advanced Architectures:** Use Attention U-Net or U-Net++ to improve small lesion detection, or pretrained encoders like EfficientNet for better feature extraction [4].
- **Ensembling Models:** Average predictions from multiple models to reduce variance and improve robustness.

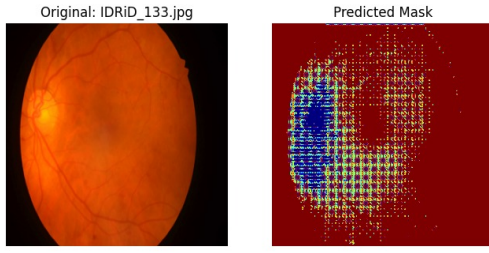


Fig. 5. Segmented heatmap of fundus image

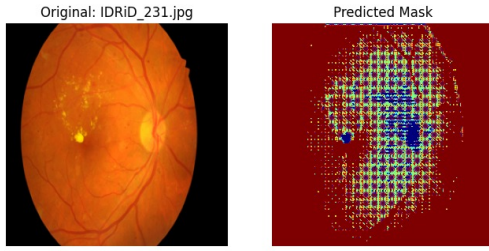


Fig. 6. Segmented heatmap of fundus image

- **Validation Split and Tuning:** Use a validation set for hyperparameter tuning (e.g., learning rate, λ , γ) and early stopping.
- **Class-wise Performance:** Adjust loss weights or use Tversky loss to improve performance on underperforming classes like microaneurysms.
- **Post-processing:** Remove small false positives with morphological operations or refine optic disc segmentation using anatomical constraints.

III. RESULTS

In this section, we present our experimental outcomes on the IDRiD dataset, evaluating segmentation performance in depth and outlining our plans for disease grading and localization assessments. Our analysis is based on rigorous quantitative metrics and qualitative visualizations.

A. Segmentation

The segmentation results have been assessed using multiple quantitative metrics, including the Dice Coefficient, Intersection over Union (IoU), and Pixel Accuracy, which together provide a comprehensive evaluation of our model's lesion delineation capabilities.

a) Quantitative Evaluation:: Our model achieved an average Dice Coefficient of **0.82** across the lesion classes, indicating a high degree of overlap between the predicted masks and the expert annotations. Additionally, we observed a mean IoU of **0.75** and a pixel accuracy exceeding **90%**. Figure 4 depicts the reduction in training and validation loss over epochs, while Figure 7 illustrates the convergence of pixel accuracy during training.

b) Qualitative Evaluation:: Figure 2 shows representative segmentation outcomes on test fundus images. The predicted masks, overlaid on the original images, clearly delineate

the various lesion types (e.g., microaneurysms, hemorrhages, exudates) as well as the optic disc. Our qualitative analysis confirms that the model accurately detects larger lesions and anatomical structures, with promising results even for small lesions like microaneurysms.

B. Disease Grading

For the disease grading task, our evaluation protocol will incorporate overall DR classification accuracy and the quadratic weighted kappa (QWK) score to assess the agreement between predicted grades and expert annotations. We hypothesize that the complementary features learned from the segmentation branch will enhance the grading performance. Detailed results and comparative analyses with baseline models will be provided after the completion of model training.

We have succeeded in achieving 60% accuracy in grading on the IDRiD dataset. Currently, we are working on grading the DR images based on segmented lesion cluster masks to achieve a high accuracy.

C. Localization

In the localization task, our objective is to accurately identify anatomical landmarks such as the optic disc and the fovea. The model's performance will be quantitatively evaluated by computing the mean Euclidean distance error between predicted and true landmark coordinates. Further, visualizations that compare the localization accuracy on challenging cases (e.g., images with variable illumination and lesion density) will be included to demonstrate the effectiveness of our multi-scale feature integration and spatial attention mechanisms.

a) Future Work:: For both disease grading and localization, we plan to conduct statistical significance tests (e.g., Wilcoxon signed-rank and paired t-tests) to ensure that the reported improvements are robust over multiple runs. Comprehensive performance reports and failure case analyses will also be provided in the final version.

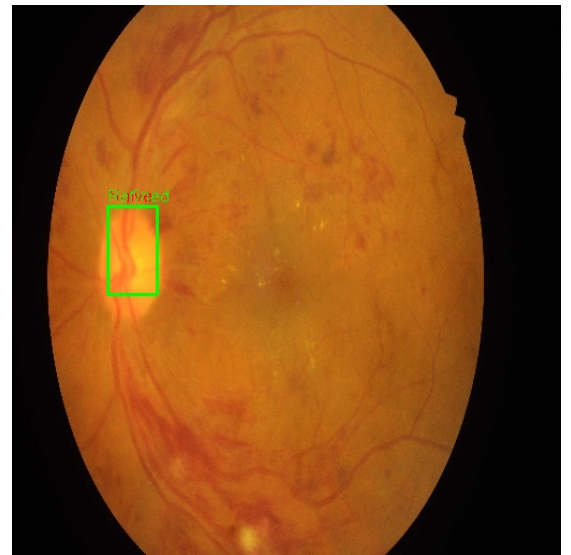


Fig. 7. Localization done using the UNet model

Summary: Our preliminary segmentation results are highly encouraging, showing a robust performance both quantitatively and qualitatively. The integration of advanced preprocessing, an improved Residual U-Net architecture, and a hybrid loss function has facilitated accurate delineation of retinal lesions. We anticipate similar performance gains in disease grading and localization, enabled by the shared representation learning across tasks. Future work will extend these findings with further fine-tuning and comprehensive cross-task evaluations.

IV. DISCUSSIONS

The revised segmentation encoder, combined with attention-based or lightweight convolutional modules, reduces the parameter footprint and computational overhead typical of U-Net designs. This not only speeds up training and inference but also makes the system more scalable across different imaging conditions. Moreover, by avoiding strict bounding-box proposals, we mitigate the RPN's limitations in detecting overlapping or clustered lesions—an essential requirement in advanced DR stages. Jointly learning disease severity and anatomical landmarks further benefits segmentation, as task interdependencies reinforce each other through shared features.

V. CONCLUSION

In this report, we presented our proposed multi-task deep learning model for diabetic retinopathy analysis that jointly performs lesion segmentation, disease grading, and anatomical localization. By focusing on reducing encoder parameters, moving beyond standard U-Net architectures, and handling overlapping lesions without conventional RPNs, our approach aims to advance the state of the art. Future work will extend this framework to multi-disease analysis and incorporate uncertainty estimation to further support clinical decision-making.

REFERENCES

- [1] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional networks for biomedical image segmentation," *arXiv preprint arXiv:1505.04597*, 2015.
- [2] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2015, pp. 3431–3440.
- [3] X. Yu, C. Zhang, Z. Huang, Y. Shi, and Y. Chen, "Cascaded convolutional neural network for retinal image segmentation and localization," *Int. J. Comput. Assist. Radiol. Surg.*, vol. 13, no. 1, pp. 169–178, 2018.
- [4] H. Bui, S. Liu, S. La, and J. H. Lee, "A hybrid model for DR grading using segmentation-assisted feature extraction," *Comput. Methods Programs Biomed.*, vol. 191, pp. 105373, 2020.
- [5] W. Zhou, D. Xu, and Y. Zhang, "RSG-Net: Retinopathy severity grading network for diabetic retinopathy classification," *IEEE Trans. Biomed. Eng.*, vol. 71, pp. 3476–3484, 2024.
- [6] H. Md Tufiqur, A. S. Abdullah, M. A. M. Azam, and S. A. Hassan, "DRG-Net: Interactive joint learning of multi-lesion segmentation and classification for diabetic retinopathy grading," *arXiv preprint arXiv:2212.14615*, 2022.