

Dear Editor and Reviewers,

Thank you for your hard work on our manuscript “DualSeg: Unified Multi-Scale Framework With Dual-Stage Encoder For Glomerular Segmentation” (ID: [JBHI-05124-2025]). We have carefully considered all your comments and made changes to the manuscript's content. In this revised version, we have significantly expanded our experimental validation by including new state-of-the-art baselines (e.g., InceptionNeXt, U-Mamba), adding a new external validation dataset (KPMP), and visualizing Effective Receptive Fields (ERF) to clarify our architectural novelty. Please note that as the manuscript was extensively condensed (from 17 to 14 pages) to meet JBHI page limits, we have restricted highlighting to substantive changes. Our responses to the comments are provided below.

## Response to reviewer #1:

### COMMENTS TO THE AUTHOR(S)

#### 1. Experimental Validation and Efficiency Claims

The paper's central efficiency claim—60% computational reduction via VRWKV—remains entirely unsubstantiated. A core contribution built on efficiency must provide comprehensive benchmarks including FLOPs, memory consumption, and inference latency across varying input sizes. The absence of these fundamental metrics suggests either inadequate experimental rigor or awareness that actual gains may not support the claimed advantages. This gap is particularly damaging given that computational efficiency differentiates DualSeg from existing methods.

**Response:** Thank you for your advice. We have addressed the concerns regarding computational benchmarking through the following revisions, as detailed below:

- **Text Revision and Clarification:** We have revised the text in **Section I. Introduction** to remove the claim regarding a "60% reduction in computational overhead." This figure was originally cited from the VRWKV framework literature (Reference [23] *Y. Duan, W. Wang, Z. Chen, X. Zhu, L. Lu, T. Lu, Y. Qiao, H. Li, J. Dai, and W. Wang, "Vision-rwkv: Efficient and scalable visual perception with rwkv-like architectures," arXiv preprint arXiv:2403.02308, 2024.*) to illustrate general efficiency, and this statement has been deleted.
- **Expanded Benchmarking Experiments:** To provide a concrete evaluation of computational efficiency, we have introduced a new comparative analysis in the revised manuscript. Specifically, we have added a **FLOPs (Floating Point Operations) vs. Dice Score comparison** across different models (now illustrated in the bottom panel of **Figure 1 and Figure 3**). This metric provides a standard and objective measure of the computational complexity and the performance-efficiency trade-off of our proposed framework.

The revisions can be found in Section I. Introduction, Figure 1 and 3:

#### Page 2, Section I. Introduction:

Conversely, Vision Transformers (ViTs) address the second challenge (spatial heterogeneity) by leveraging self-attention for global context [20], [21], as shown in Fig. 1(II). ViTs typically outperform UNets in maintaining structural continuity amidst fibrosis-induced fragmentation [21]. Yet, the quadratic complexity of self-attention limits their utility in high-resolution histology [10]. To mitigate this, efficient alternatives like Wave-MLP [22] (preserving structure via spatial wise convolutions) and VRWKV [23] (linear-time recurrent kernels) have emerged. Similarly, VM-UNet [24] introduced State Space Models (SSM) to harmonize efficiency and performance.

...

To answer this, we propose DualSeg, a novel hybrid framework synergizing Wave Vision and VRWKV within a pyramid structure (Fig. 2(f)). Specifically, our architecture employs a dual-stage encoder: early-stage Wave-Swin blocks perform hierarchical local feature extraction to resolve texture discriminability, while later-stage VRWKV blocks model long-range dependencies via linear attention to address spatial heterogeneity. This design combines the texture sensitivity of CNNs, the generalization of MLPs, and the scalability of VRWKV. By bridging local and global processing, DualSeg achieves robust multi-scale mapping of morphological priors. As consistently demonstrated across **Fig. 1 and Fig. 3**, DualSeg synergizes a global, clean ERF (**Fig. 1(IV)**) with SOTA segmentation performance while maintaining a minimal computational footprint (3.09 G FLOPs), establishing a superior accuracy-efficiency balance compared to resource-intensive baselines.

## Page 1, Figure 1:

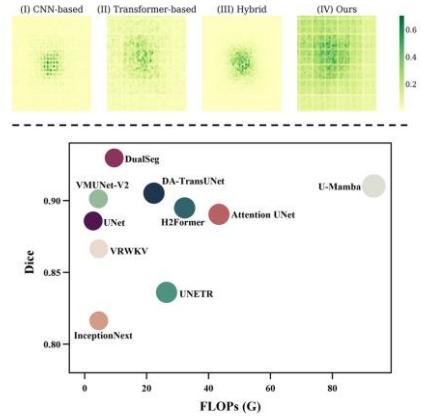


Fig. 1: **Top:** Visualization of the *Effective Receptive Fields* (ERF) for different architectures. CNNs (I) focus locally, while Transformers (II) capture global but noisy patterns. Our method (IV) achieves a clean, global ERF. **Bottom:** Performance vs. FLOPs comparison. DualSeg (top-left) achieves the optimal trade-off between segmentation accuracy and computational efficiency.

## Page 2, Figure 3:

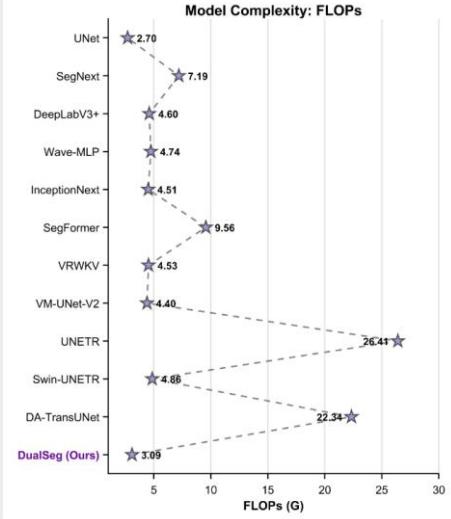


Fig. 3: Benchmarking of computational complexity (G FLOPs). Our DualSeg achieves an optimal balance with only 3.09 G FLOPs. For better visualization of the fine-grained differences among efficient models, outliers with computational costs > 30 G FLOPs are excluded from this plot.

## 2. Dataset Limitations and Clinical Generalizability

The evaluation's restriction to PAS-stained specimens fundamentally limits clinical relevance. Real-world pathology predominantly uses H&E staining with supplementary protocols (Masson's trichrome, Jones silver) for specific diagnoses. Each staining method reveals distinct tissue characteristics—algorithms optimized for PAS often fail catastrophically on H&E due to different contrast patterns and feature visibility. The absence of cross-institutional validation or multi-protocol testing renders clinical applicability claims premature.

**Response:** Thank you for pointing that out. We have revised the manuscript to strengthen the evaluation of generalizability and address staining limitations, as detailed below:

- **Expanded Validation Experiments:** We have incorporated a new external dataset—the **KPMP (Kidney Precision Medicine Project) dataset** (now detailed in **Section IV. A. Dataset**). This allowed us to perform an extensive validation encompassing: **Cross-Institutional Validation:** Testing on data from completely different centers without any fine-tuning (zero-shot inference). The results demonstrate that even when restricted to PAS staining, DualSeg maintains high robustness against significant biological and technical variations across different cohorts. And the results are illustrated in **Table III and Figure 9**.
- **Text Revision on Staining Limitations:** We have added a discussion in **Section VI. D. Limitations and Future Work**, regarding the scarcity of synchronized multi-stain (e.g., H&E, Masson) annotations in current public datasets. The revised text explicitly states that while DualSeg demonstrates superior performance on PAS-stained images, its clinical deployment across diverse staining protocols remains a subject for future validation as relevant data becomes available.

The revisions can be found in Section IV. A. Dataset, Section VI. D. Limitations and Future Work, Table III and Figure 8:

### Page 6, Section IV. A. Datasets:

Dataset III: Human Glomeruli (KPMP). To assess cross species generalization, we retrieved a second human dataset from the Kidney Precision Medicine Project (KPMP) Atlas Repository [50]. Four PAS-stained SVS format WSIs (avg. resolution  $84,000 \times 50,000$ ) were selected with corresponding masks. To rigorously validate generalization, models trained solely on the mouse KPIs dataset were directly applied to this human dataset without retraining. For preprocessing consistency, KPMP WSIs were partitioned into  $2,048 \times 2,048$  patches.

### Page 13, Section VI. D. Limitations and Future Work:

While DualSeg demonstrates superior performance in glomerular segmentation, three primary limitations remain to be addressed in future iterations. First, the model's generalization to rare pathological subtypes, such as global glomerulosclerosis, is currently constrained by data scarcity. We plan to mitigate this by employing domain adaptation and transfer learning techniques to enhance feature robustness for these underrepresented classes. Second, our current validation is restricted to PAS-stained images. To ensure broad clinical applicability, we will extend our evaluation to include Hematoxylin-Eosin (HE) stained datasets, verifying the model's adaptability to varying histological protocols. Finally, despite the linear complexity of the VRWKV block, the computational overhead for gigapixel WSI processing remains significant. Future work will focus on model quantization and lightweight optimization to facilitate deployment on resource-constrained platforms and edge devices.

### Page 7, Table III:

TABLE III: CROSS-DATASET INFERENCE PERFORMANCE COMPARISON FOR GLOMERULAR SEGMENTATION ON THE KMPM DATASET USING 5-FOLD MOUSE-TRAINED MODELS WITH RESPECT TO EXISTING METHODS

Models	1			2			3			4			AVG		
	mDSC $\dagger$	HD95 $\dagger$	IoU $\dagger$	mDSC $\dagger$	HD95 $\dagger$	IoU $\dagger$	mDSC $\dagger$	HD95 $\dagger$	IoU $\dagger$	mDSC $\dagger$	HD95 $\dagger$	IoU $\dagger$	mDSC $\dagger$	HD95 $\dagger$	IoU $\dagger$
U-Net [7]	0.3936 ±0.0045	279.4223 ±2.0000	0.5095 ±0.0045	0.2662 ±0.0042	471.0985 ±2.1642	0.2560 ±0.0042	0.3189 ±0.0047	119.5596 ±1.8510	0.4973 ±0.0047	0.5387 ±0.0046	152.3241 ±2.0442	0.5090 ±0.0046	0.4836 ±0.0044	196.2390 ±2.0462	0.4588 ±0.0046
Attention U-Net [41]	0.7165 ±0.0038	165.1443 ±2.0000	0.6875 ±0.0038	0.4762 ±0.0046	274.0519 ±2.0000	0.4540 ±0.0046	0.4831 ±0.0043	54.8253 ±1.7304	0.8230 ±0.0043	0.6953 ±0.0043	90.7108 ±2.0000	0.6559 ±0.0044	0.7056 ±0.0044	110.6999 ±1.9678	0.6736 ±0.0041
SegNest [26]	0.7026 ±0.0036	164.3759 ±2.0000	0.7656 ±0.0036	0.5203 ±0.0046	252.5286 ±2.0000	0.6078 ±0.0046	0.5888 ±0.0043	73.0842 ±1.7304	0.7877 ±0.0043	0.7987 ±0.0043	83.8105 ±2.0000	0.6948 ±0.0043	0.812102 ±0.0040	0.7100 ±0.0040	0.7100 ±0.0040
DeepLabV3+ [34]	0.7846 ±0.0038	100.0552 ±1.9003	0.7640 ±0.0038	0.6030 ±0.0043	182.8433 ±2.0000	0.5615 ±0.0043	0.5845 ±0.0043	67.7374 ±1.6939	0.8333 ±0.0043	0.6952 ±0.0043	116.7432 ±1.6939	0.6617 ±0.0043	0.7317 ±0.0043	112.6154 ±1.9438	0.7018 ±0.0040
Wave-MLP [22]	0.8152 ±0.0040	149.4421 ±2.0000	0.7899 ±0.0038	0.5158 ±0.0043	299.4355 ±2.0000	0.4929 ±0.0043	0.5095 ±0.0043	49.6731 ±1.6939	0.8253 ±0.0043	0.7675 ±0.0043	92.2304 ±1.6939	0.7342 ±0.0043	0.7646 ±0.0043	106.4041 ±1.9438	0.7353 ±0.0040
InceptionNext [42]	0.6564 ±0.0044	236.3450 ±2.0000	0.6352 ±0.0044	0.4779 ±0.0046	272.6340 ±2.0000	0.4672 ±0.0046	0.7130 ±0.0046	77.7567 ±1.7304	0.6037 ±0.0046	0.4301 ±0.0046	312.1570 ±2.0000	0.3875 ±0.0046	0.5279 ±0.0046	294.8018 ±1.9438	0.4967 ±0.0046
SegFormer [11]	0.8037 ±0.0036	125.7612 ±1.8166	0.7814 ±0.0037	0.6151 ±0.0044	206.3655 ±2.0000	0.5815 ±0.0044	0.5897 ±0.0044	47.1347 ±1.0667	0.8353 ±0.0044	0.8285 ±0.0044	49.0933 ±1.1678	0.7986 ±0.0044	0.8082 ±0.0043	70.7543 ±1.4692	0.7802 ±0.0043
VRWKV [23]	0.7012 ±0.0036	143.0609 ±2.0000	0.7662 ±0.0035	0.5112 ±0.0046	271.1789 ±2.0000	0.4851 ±0.0046	0.5056 ±0.0046	62.1066 ±1.0115	0.7850 ±0.0046	0.7580 ±0.0046	65.9837 ±2.0000	0.7237 ±0.0046	0.7480 ±0.0046	80.7188 ±1.9438	0.7188 ±0.0040
VM-UNET-V2 [40]	0.6521 ±0.0043	181.1196 ±1.9553	0.6233 ±0.0044	0.4833 ±0.0046	332.8376 ±2.0000	0.4575 ±0.0046	0.7880 ±0.0046	104.5297 ±1.5290	0.7620 ±0.0046	0.5410 ±0.0046	232.7710 ±2.4792	0.6062 ±0.0043	0.6028 ±0.0043	216.0202 ±2.3843	0.5664 ±0.0044
UNETR [9]	0.6199 ±0.0046	250.5464 ±2.0000	0.6050 ±0.0046	0.4500 ±0.0046	391.7400 ±2.0000	0.4416 ±0.0046	0.6436 ±0.0046	158.4295 ±1.6349	0.6349 ±0.0046	0.4890 ±0.0046	333.3864 ±2.0000	0.4416 ±0.0046	0.5369 ±0.0046	308.9828 ±2.0276	0.5193 ±0.0047
Swin UNETR [43]	0.5108 ±0.0047	207.3004 ±2.0000	0.5261 ±0.0046	0.3588 ±0.0046	432.1604 ±2.0000	0.3240 ±0.0046	0.7057 ±0.0046	85.0000 ±1.7000	0.5600 ±0.0046	0.4699 ±0.0046	131.7901 ±2.0001	0.4921 ±0.0046	0.5320 ±0.0046	152.2606 ±1.9438	0.5320 ±0.0046
DA-TransU-Net [44]	0.7783 ±0.0038	97.7011 ±1.7154	0.7563 ±0.0039	0.5007 ±0.0045	166.4438 ±2.1574	0.5010 ±0.0045	0.8310 ±0.0045	63.6652 ±1.4344	0.8566 ±0.0045	0.8143 ±0.0045	81.5234 ±1.4344	0.7594 ±0.0045	0.7780 ±0.0045	87.6505 ±1.6488	0.7489 ±0.0038
H2Former [10]	0.7490 ±0.0040	148.2938 ±2.0000	0.7220 ±0.0040	0.5027 ±0.0046	219.4662 ±2.0000	0.5223 ±0.0046	0.8095 ±0.0046	55.2702 ±1.5290	0.7859 ±0.0046	0.6369 ±0.0046	169.6482 ±2.0000	0.5934 ±0.0046	0.6815 ±0.0046	154.2361 ±1.9438	0.6460 ±0.0042
U-mamba [25]	0.6549 ±0.0043	126.1347 ±1.7237	0.6214 ±0.0043	0.5967 ±0.0045	204.4849 ±1.6164	0.5729 ±0.0045	0.7500 ±0.0045	1.1344 ±1.5002	0.7282 ±0.0045	0.5589 ±0.0045	282.8070 ±2.3296	0.5875 ±0.0045	0.6101 ±0.0045	53.4040 ±2.1725	0.5440 ±0.0043
DualSegs (Ours)	<b>0.8251</b> ±0.0034	<b>81.6201</b> ±1.5314	<b>0.8022</b> ±0.0035	<b>0.5888</b> ±0.0045	<b>235.8327</b> ±2.6482	<b>0.5573</b> ±0.0045	<b>0.8848</b> ±0.0028	<b>21.0608</b> ±0.2919	<b>0.8630</b> ±0.0029	<b>0.8393</b> ±0.0032	<b>57.6049</b> ±1.3786	<b>0.8123</b> ±0.0032	<b>0.8195</b> ±0.0034	<b>69.6369</b> ±1.5420	<b>0.7938</b> ±0.0035

Page 11, Figure 9:

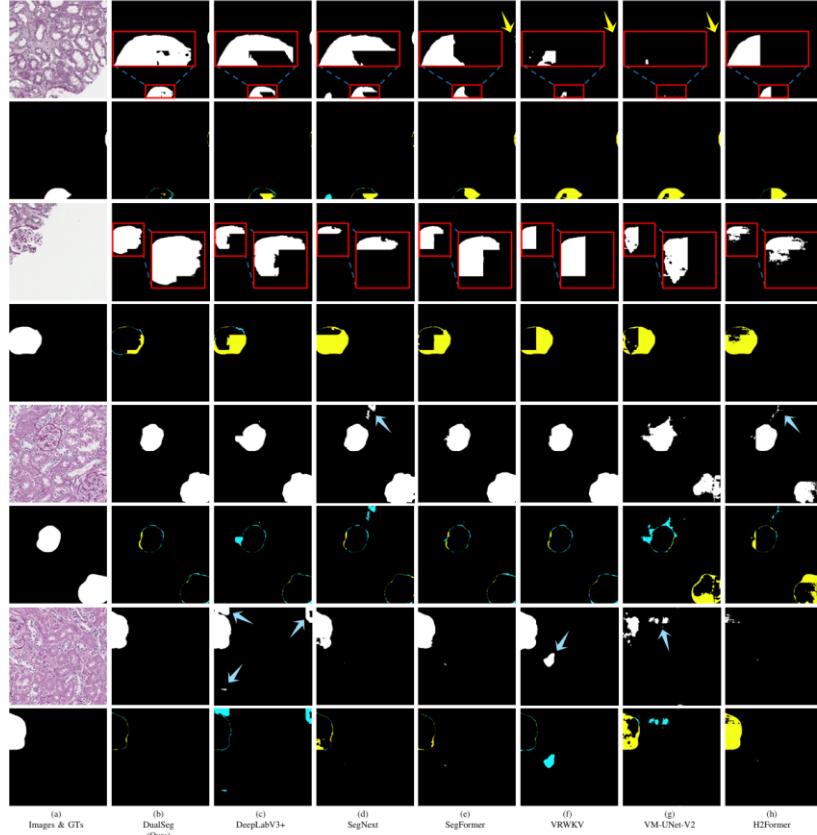


Fig. 8: Visual comparison on the held-out KMPM test set. **Odd rows** display original images and inference masks; **even rows** show GT and error maps, where **yellow** and **green** denote under-segmentation and over-segmentation, respectively. Red boxes and matching arrows highlight magnified details and segmentation errors. **DualSeg** exhibits unprecedented stability when handling cross-center and cross-species data.

### 3. Incomplete Architectural Comparisons

The baseline selection reveals a critical methodological flaw through systematic omission of contemporary efficient architectures. InceptionNeXt, which achieves transformer-comparable performance with CNN efficiency, directly challenges DualSeg's value proposition yet remains unexamined. Similarly, Mamba-based segmentation models (VM-UNet, U-Mamba) already address the linear complexity challenge DualSeg claims to solve. These omissions appear deliberate rather than oversight, suggesting the authors recognize these comparisons might undermine their architectural superiority claims. Without these essential benchmarks, the true contribution remains indeterminate.

**Response:** Thank you for pointing that out. We have addressed the perceived "methodological gap" by conducting a comprehensive re-evaluation. In the revised manuscript, we have integrated contemporary efficient architectures, including InceptionNeXt, VM-UNet, and U-Mamba, as primary baselines. These models are now formally described in **Section IV. B. Baselines**. The quantitative and qualitative results of these comparisons are presented in **Section V. A. Glomeruli Segmentation Results** and **Section V. C. Visualization Results**, detailed in **Tables I, II, and III**, and visualized in **Figure 8**. A detailed analysis and comparative discussion regarding our model's superiority over these state-of-the-art baselines are provided in **Section VI. A. Glomeruli Segmentation Results**. Our findings consistently demonstrate that DualSeg achieves superior segmentation accuracy while maintaining a more favorable performance-to-complexity ratio.

The revisions can be found in Section IV. B. Baselines, Section VI. A. Discussion, Section V. A. Results, Tables I-III and Figure 9:

#### Page 7, Section IV. B. Baselines:

We benchmark DualSeg against 14 representative methods spanning three architectural paradigms. CNN-based models include canonical baselines like U-Net [7] and Attention U-Net [41], the receptive-field-enhanced DeepLabV3+ [34], and efficient modern architectures such as SegNext [26], InceptionNext [42], and Wave-MLP [22]. Transformer-based models encompass SegFormer [11] for hierarchical encoding, VRWKV [23] utilizing linear recurrent operators, and the **SSM-integrated VM-UNet-V2** [40]. Finally, Hybrid models feature U-shaped Transformer variants like UNETR [9] and SwinUNETR [43], alongside advanced fusion frameworks including H2Former [10], DA-TransUNet [44], and the **Mamba-based U-Mamba** [25].

#### Page 9, Section VI. A. Glomeruli Segmentation Results:

3) KPMP Dataset (Cross-Species Inference): To assess generalizability, we evaluated models pre-trained solely on the murine KPIs dataset directly on the human KPMP dataset without additional fine-tuning (Table III). **DualSeg achieved a highly competitive average mDSC of 81.95%, surpassing SegFormer (80.82%) and significantly outperforming InceptionNext (52.79%).** In terms of boundary delineation, DualSeg achieved an HD95 of 69.64, which is approximately one third of the error recorded by VM-UNet-V2 (216.02 units). Furthermore, its IoU of 79.38% surpassed other hybrid models by over 4%, validating the model's capacity for robust cross-species and cross-center transfer learning.

#### Page 11, Section VI. A. Comparative Analysis with SOTA Methods:

DualSeg effectively addresses the limitations of existing architectural paradigms through a unified framework. **Unlike conventional CNNs (e.g., DeepLabV3+, InceptionNext) which are constrained by fixed receptive fields, DualSeg utilizes the Wave-Swin Block's dynamic propagation window (Table V).** This innovation enhances local texture discriminability, reducing HD95 by 26–37 on the KPIs dataset compared to CNN baselines and outperforming InceptionNext by 29.16% mDSC in cross-species tasks.

Furthermore, while Transformers typically excel at global context but compromise local detail due to patch flattening, DualSeg integrates a Z-Shift operator within the VRWKV block to preserve edge integrity. **This design mitigates the information loss inherent in standard Q-Shift operations (Table VI) and models spatial heterogeneity more effectively, allowing DualSeg to exceed VM-UNet-V2 by 21.67% on the external KPMP dataset.**

Finally, in contrast to existing hybrids (e.g., H2Former, DA TransUNet) that suffer from a “semantic gap” or the high computational costs of U-Mamba (Fig. 1), DualSeg employs a sequential *local-to-global* refinement strategy. This structured integration ensures precise morphological prior mapping, enabling the model to outperform H2Former by 13.68% on HuBMAP and surpass U-Mamba by 23.52% on KPMP. By dynamically adapting to morphological variability and seamlessly integrating features, DualSeg establishes a robust and efficient backbone for renal histology analysis.



Page 11, Figure 9:

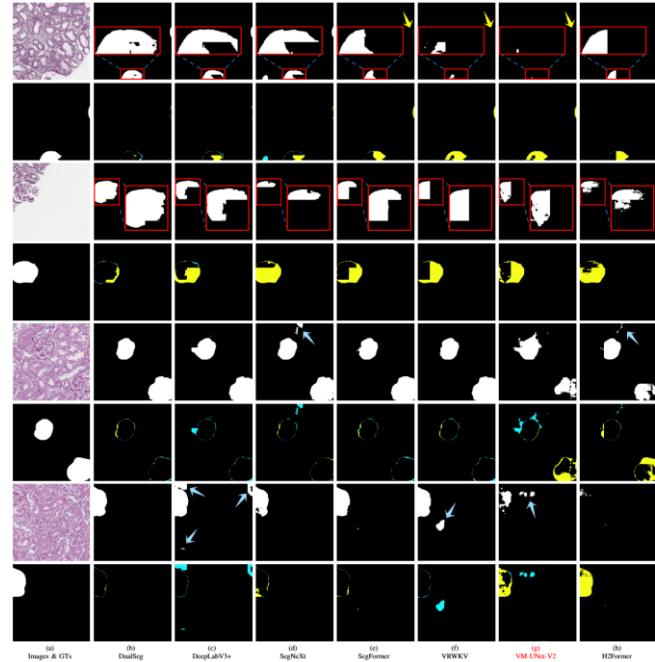


Fig. 9: Visual comparison on the held-out KPMP test set. **Odd rows** display original images and inference masks; **even rows** show GT and error maps, where **yellow** and **green** denote under-segmentation and over-segmentation, respectively. Red boxes and matching arrows highlight magnified details and segmentation errors. **DualSeg** exhibits unprecedented stability when handling cross-center and cross-species data.

#### 4. Clinical Deployment Analysis

Despite positioning DualSeg for practical application, the manuscript provides no deployment feasibility analysis. Clinical environments impose strict constraints: limited GPU memory (often 8GB), CPU-only workstations in many facilities. The reported 2.73% mDSC improvement lacks clinical context—pathologists require understanding of whether such margins affect diagnostic confidence, inter-observer agreement, or treatment decisions. Without this translation from metrics to clinical impact, the practical value remains speculative.

**Response:** Thank you for pointing that out. We have revised the manuscript to demonstrate the statistical robustness and practical feasibility of our results, as detailed below:

- **Text Revision on Deployment Feasibility:** We have added a comprehensive discussion in **Section VI. D. Limitations and Future Work** regarding the practical deployment of DualSeg. This new section outlines the roadmap for future optimizations, including lightweight design to transition our architectural innovations into clinical tools.
- **Expanded Statistical Validation:** To substantiate the 2.73% mDSC improvement, we have enriched **Section VI. C. Clinical Relevance** with quantitative evidence. We performed *t*-tests across three diverse datasets (KPIs, HuBMAP, and KPMP); as illustrated in **Figure 11**, the results confirm that DualSeg’s performance gains are statistically significant ( $p < 0.001$ ), ensuring the improvements are robust rather than marginal.
- **Enhanced Qualitative Evaluation:** To supplement the quantitative metrics, we have incorporated a granular spatial audit in **Figures 7 and 9**. By utilizing red bounding boxes for local magnification and colored arrows (yellow for under-segmentation; green for over-segmentation), we explicitly demonstrate DualSeg’s superior ability to maintain structural continuity and reduce critical diagnostic omissions in complex anatomical structures.

The revisions can be found in Section VI. D. Limitations and Future Work, Section VI. C. Clinical Relevance, Figures 7, 9 and 11:

##### **Page 13, Section VI. D. Limitations and Future Work:**

While DualSeg demonstrates superior performance in glomerular segmentation, three primary limitations remain to be addressed in future iterations. First, the model’s generalization to rare pathological subtypes, such as global glomerulosclerosis, is currently constrained by data scarcity. We plan to mitigate this by employing domain adaptation and transfer learning techniques to enhance feature robustness for these underrepresented classes. Second, our current validation is restricted to PAS-stained images. To ensure broad clinical applicability, we will extend our evaluation to include Hematoxylin-Eosin (HE) stained datasets, verifying the model’s adaptability to varying histological protocols. **Finally, despite the linear complexity of the VRWKV block, the computational overhead for gigapixel WSI processing remains significant. Future work will focus on model quantization and lightweight optimization to facilitate deployment on resource-constrained platforms and edge devices.**

##### **Page 13, Section VI. C. Clinical Relevance:**

DualSeg exhibits statistically significant superiority ( $p < 0.05–0.001$ ; Fig. 11) and exceptional reproducibility, evidenced by a minimal standard deviation (0.0040) on the HuBMAP dataset. Its ability to accurately resolve diverse morphologies—ranging from mild hypertrophy to severe fragmentation—enables the precise quantification of pathological biomarkers like sclerosis and fibrosis. Furthermore, the model’s robust performance on the cross-species KPMP dataset supports standardized CKD monitoring. By mitigating inter-observer variability and reducing manual annotation burdens, DualSeg provides a scalable solution for multi-center clinical trials and routine diagnostic workflows.

Page 8, Figure 7:

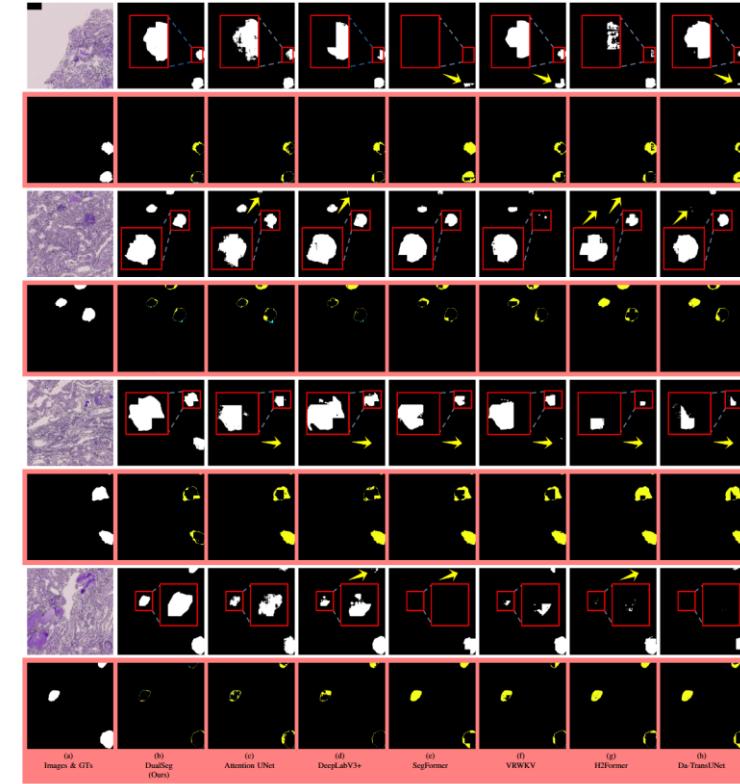


Fig. 7: Visual comparison on the KPIs test set. **Odd rows** display original images and inference masks; **even rows** show GT and error maps, where **yellow** and **green** denote under-segmentation and over-segmentation, respectively. Red boxes and matching arrows highlight magnified details and specific segmentation errors. **DualSeg** demonstrates superior performance in handling heterogeneous glomeruli with intricate boundaries.

Page 11, Figure 9:

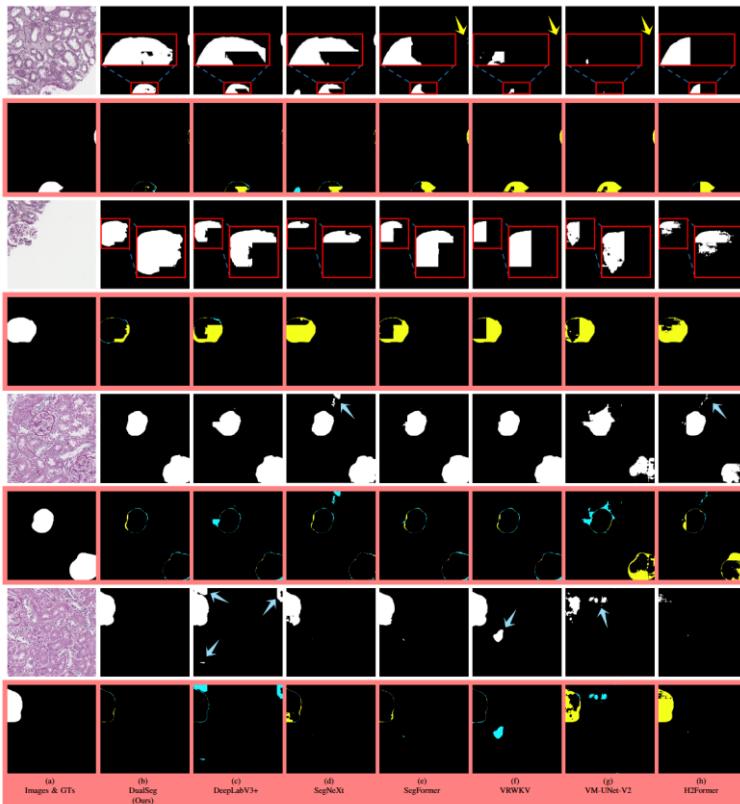


Fig. 9: Visual comparison on the held-out KPMP test set. **Odd rows** display original images and inference masks; **even rows** show GT and error maps, where **yellow** and **green** denote under-segmentation and over-segmentation, respectively. Red boxes and matching arrows highlight magnified details and segmentation errors. **DualSeg** exhibits unprecedented stability when handling cross-center and cross-species data.

Page 12, Figure 11:

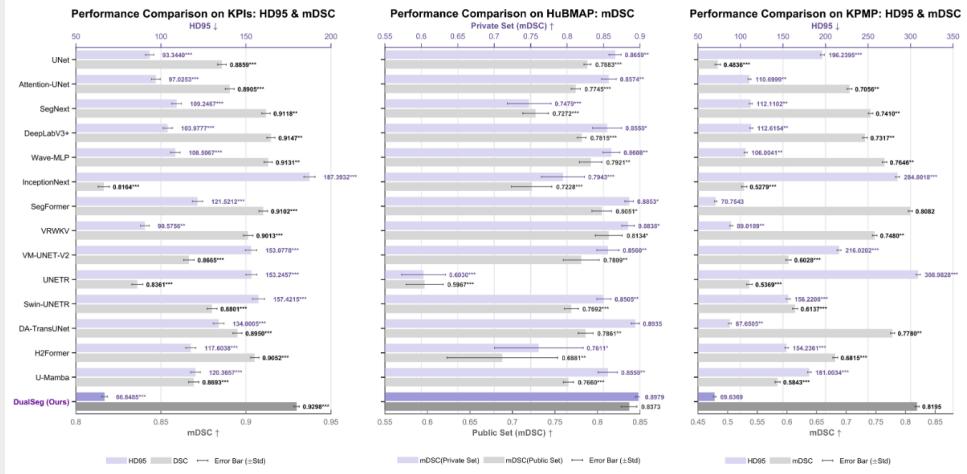


Fig. 11: Performance comparison between our DualSeg model and 14 baseline methods across three datasets. The comparison metrics include mDSC and HD95 for the KPIs and KPMP datasets (left and right panels, respectively), and mDSC for the HuBMAP dataset (middle panel). The error bars represent  $\pm$  standard deviation. Statistical significance was assessed using paired *t*-tests, with levels indicated by asterisks: \* $p < 0.05$ , \*\* $p < 0.01$ , \*\*\* $p < 0.001$ .

## 5. Architectural Design Justification

The dual-stage encoder's sequential arrangement (Wave-Swin → VRWKV) lacks theoretical foundation. The manuscript presents this as optimal without exploring alternative configurations (parallel processing, reversed ordering, or hybrid fusion strategies). The dynamic window sizing, acknowledged as "manually defined based on empirical observations," reveals methodological weakness—critical design parameters derived through trial-and-error rather than principled analysis. This empirical approach, while sometimes necessary, requires thorough justification absent here.

**Response:** Thank you for your valuable comments. We have addressed the concerns regarding the model's practical utility by adding a limitation analysis on deployability and a dedicated section on clinical significance, as detailed below:

- **Ablation of Architectural Ordering:** To substantiate the theoretical superiority of the proposed sequential arrangement (Wave-Swin→VRWKV), we have conducted additional ablation experiments evaluating alternative configurations, including reversed ordering (VRWKV→Wave-Swin) and Attention-Wave hybrids. The comparative results (detailed in the **Section V. B. Ablation Studies** and **Table IV**) demonstrate that our current design optimizes the transition from local spatial feature extraction to global linear complexity modeling, yielding the highest segmentation accuracy.
- **Principled Justification for Dynamic Window Sizing:** We clarify that the selection of the window candidate set  $S = \{7, 11, 15, \dots\}$  is not the result of unprincipled trial-and-error, but is rooted in domain-validated baselines. We have significantly expanded **Section III. A. 2. Dynamic Swin Mechanism** to detail this rationale.

The revisions can be found in Section III. A. Wave-Swin Block, Section V. B. Ablation Studies and Table IV:

### Page 4, Section III. A. Wave-Swin Block:

First, previous studies on Wave-MLP have empirically demonstrated that windows smaller than 7 lack the generality necessary to capture spatial dependencies in medical images [22]; meanwhile, anchor sizes 7 and 11 align with kernel sizes employed in SOTA encoders like SegNeXt [26]. Second, the average glomerular bounding box in our murine dataset measures approximately 154px [33]. After the 4x and 8x downsampling stages, this dimension reduces to roughly 38px and 19px, respectively. Accordingly, selecting a maximum window of 15 (instead of 21) prevents the network from integrating extraneous background noise while ensuring full coverage of the target glomerular structure.

### Page 9, Section V. B. Ablation Studies:

1) Effect of Dual-Stage Encoder: Table IV contrasts single stage architectures with the proposed dual-stage design. Single-stage variants employing only Wave-Swin Blocks, SegFormer-style self-attention, or VRWKV Blocks achieved average mDSCs of 90.36%, 90.99%, and 91.08%, respectively. The integrated Dual-Stage (Wave-VRWKV) architecture outperformed all single-stage counterparts with an average mDSC of 92.98%. We further investigated the impact of module sequencing. Reversing the feature extraction order (placing attention mechanisms before wave blocks) resulted in a significant performance decrease, lowering the average mDSC to 85.78%(Attention-Wave) and 91.63%(VRWKV-Wave). These findings corroborate the critical role of the proposed local-to-global refinement strategy.

**Page 9, Table IV:**

TABLE IV: ABLATION STUDY OF MAJOR COMPONENTS ON THE TEST SET OF THE KPIs DATASET

Stage	Models	Layers			mDSC↑				
		Wave	Attention	VRWKV	DN	NEP25	Normal	S/Nx	Avg
Sole-Stage	Wave [22]	✓	-	-	0.9165	0.9236	0.9322	0.8068	0.9036
	Attention [11]	-	✓	-	0.9223	0.9123	0.9249	0.8603	0.9099
	VRWKV [23]	-	-	✓	0.9209	0.9117	0.9268	0.8613	0.9108
Dual-Stage	Attention-Wave(Ours)	✓	✓	-	0.9192	0.9086	0.9241	0.8275	0.9019
	VRWKV-Wave(Ours)	✓	-	✓	0.8578	0.8168	0.8694	0.7011	0.8271
	Wave-Attention(Ours)	✓	✓	-	0.9267	0.9174	0.9334	0.8678	0.9171
	Wave-VRWKV(Ours)	✓	-	✓	0.9603	0.9349	0.9187	0.9007	0.9298

## Minor Concerns

**Ablation scope:** Component analysis limited to mDSC ignores computational overhead—does each module justify its complexity?

- **Response:** We have addressed the concern regarding computational overhead by incorporating a system-level efficiency analysis in **Figure 1 (Bottom) and 3**, which demonstrates that DualSeg achieves an optimal trade-off between segmentation accuracy and FLOPs compared to 14 baseline methods. While our ablation studies (**Tables IV-VI**) prioritize the synergy of the dual-stage architecture, the selection of the lightweight HamDecoder using Non-negative Matrix Factorization (NMF) further justifies our design by enhancing multi-scale fusion without significant complexity. This strategic balance ensures that the proposed local-to-global refinement remains computationally viable for high-resolution histopathology.

The supporting data for these revisions can be found in Figure 1, 3 (pertaining to Comment 1), Table IV (pertaining to Comment 5), and Tables V and VI, which are provided below:

### Page 9, Tables V and VI:

TABLE V: ABLATION STUDY OF THE PROPAGATION WINDOW ON THE TEST SET OF THE KPIs DATASET

Models	Propagation Window Size	mDSC↑				
		DN	NEP52	Normal	S/6Nx	AVG
Wave-MLP [22]	7	0.9299	0.9189	0.9361	0.8375	0.9131
	11	0.9228	0.9175	0.9366	0.8430	0.9075
	15	0.9227	0.9189	0.9144	0.8320	0.9012
	7-15	0.9365	0.9236	0.9322	0.8645	0.9146
DualSeg(Ours)	7	0.9544	0.9332	0.9245	0.8737	0.9205
	11	0.9571	0.9123	0.9224	0.8571	0.9132
	15	0.9450	0.9122	0.9111	0.8480	0.9112
	7-15	0.9603	0.9349	0.9187	0.9007	0.9298

TABLE VI: ABLATION STUDY OF THE SHIFT MODE ON THE TEST SET OF THE KPIs DATASET

Models	Shift Mode	mDSC↑				
		DN	NEP52	Normal	S/6Nx	AVG
VRWKV [23]	<i>Q-Shift</i>	0.9330	0.9232	0.9212	0.8738	0.9013
	<i>Z-Shift</i>	0.9344	0.9377	0.9255	0.8840	0.9108
DualSeg(Ours)	<i>Q-Shift</i>	0.9554	0.9220	0.9166	0.8902	0.9177
	<i>Z-Shift</i>	0.9603	0.9349	0.9187	0.9007	0.9298

**Presentation imbalance:** Technical density overshadows clinical motivation, limiting accessibility.

- **Response:** Thank you for the constructive suggestions. To better balance technical density with clinical utility, we have enriched **Section VI. C. Clinical Relevance** with a statistical analysis of diagnostic impact. As illustrated in the newly added **Fig. 11**, DualSeg demonstrates statistically significant superiority ( $p < 0.05$  to  $p < 0.001$ ) across three diverse datasets, confirming that the mDSC improvements are robust and clinically meaningful. These performance gains directly translate to a reduction in diagnostic omissions and more precise quantification of biomarkers like glomerulosclerosis, thereby effectively reducing the manual annotation burden for pathologists.

The supporting data for these revisions can be found in Figure 11 (pertaining to Comment 4) and Section VI. C. Clinical Relevance, which is provided as follows:

### Page 13, Section VI. C. Clinical Relevance:

DualSeg exhibits statistically significant superiority ( $p < 0.05$ – $0.001$ ; Fig. 11) and exceptional reproducibility, evidenced by a minimal standard deviation (0.0040) on the HuBMAP dataset. Its ability to accurately resolve diverse morphologies—ranging from mild hypertrophy to severe fragmentation—enables the precise quantification of pathological biomarkers like sclerosis and fibrosis. Furthermore, the model’s robust performance on the cross-species KPMP dataset supports standardized CKD monitoring. By mitigating inter-observer variability and reducing manual annotation burdens, DualSeg provides a scalable solution for multi-center clinical trials and routine diagnostic workflows.

**Visualization gaps:** Comparison figures lack error maps or uncertainty quantification essential for understanding performance differences.

- **Response:** Thank you for the advice. To address the gap in performance interpretation, we have incorporated **Error Maps** into the comparative visualizations in **Figures 7 and 9**. These maps provide a spatial quantification of segmentation uncertainty and errors, allowing for a more nuanced understanding of where the model excels and where its limitations lie compared to baseline methods.

The supporting data for these revisions can be found in Figure 7 and 9 (pertaining to Comment 4).

## Response to reviewer #2:

### COMMENTS TO THE AUTHOR(S)

This manuscript, entitled “DualSeg: Unified multi-scale framework with dual-stage encoder for glomerular segmentation,” presents a dual-stage hybrid segmentation model that combines convolutional and recurrent attention mechanisms (Wave-Swin and VRWKV) to improve glomerular segmentation performance in kidney histopathology images. The topic is timely and relevant to renal pathology and computational histology. The model demonstrates good quantitative results on both mouse and human datasets.

However, despite these merits, the work still contains several limitations in both methodology and presentation, and a major revision is required before it can be considered for publication.

- 1) **the claimed novelty is not entirely convincing.** The proposed framework, while integrating CNN and VRWKV modules, appears conceptually similar to previously published hybrid architectures such as TransUNet, H2Former, and DA-TransUNet. The paper does not provide sufficient theoretical or empirical evidence to show that DualSeg fundamentally differs from these approaches. The authors should clearly articulate the unique contribution of their dual-stage design, beyond incremental improvements or architectural recombination. It would be helpful to include visual or quantitative analysis (e.g., feature map visualization, attention distribution comparison) to demonstrate how the proposed design contributes to performance beyond existing hybrid models.

**Response:** Thank you for your advice. We have addressed the concern regarding architectural novelty by clarifying the fundamental differences between DualSeg and existing hybrid designs like TransUNet and H2Former. Unlike traditional U-shaped cascades, DualSeg employs a hierarchical pyramid structure that sequentially integrates local wave-based texture refinement and global linear-complexity modeling (VRWKV). To provide empirical evidence of this advantage, we have included a visualization of the **Effective Receptive Field (ERF)** in **Fig. 1 (Top)**. The comparison demonstrates that while prior hybrid architectures often exhibit restricted or noisy receptive fields (**Fig. 1-III**), our dual-stage design achieves a distinctively "clean and global" ERF (**Fig. 1-IV**). This visualization confirms that DualSeg uniquely bridges the gap between local precision and global structural continuity, facilitating more robust feature extraction than previous architectural recombinations.

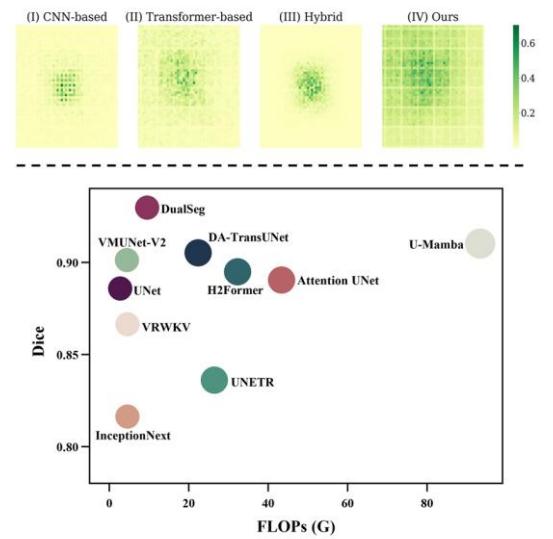
The revisions can be found in Section I. Introduction and Figure 1:

#### Page 2, Section I. Introduction:

Beyond standalone architectures, synergistic hybrid frame works have been explored. As illustrated in Fig. 2(a-d), most prior hybrids (e.g., TransUNet [8], H2Former [10] and U mamba [25]) adopt U-shaped paradigms. However, these often suffer from limitations: TransUNet compromises multi-scale capture; H2Former's shallow integration underutilizes ViTs; and their ERFs often remain suboptimal (Fig. 1(III)). Distinct from U-shaped models, pyramid architectures like SegFormer [11] (Fig. 2(e)) and SegNext [26] address the third challenge via feature fusion but rely on unidirectional extraction, weakening robustness against morphological variations. This prompts a critical inquiry: *Is it possible to integrate local and global features within a unified multi-scale framework to tackle all three segmentation challenges simultaneously?*

To answer this, we propose DualSeg, a novel hybrid framework synergizing Wave Vision and VRWKV within a pyramid structure (Fig. 2(f)). Specifically, our architecture employs a dual-stage encoder: early-stage Wave-Swin blocks perform hierarchical local feature extraction to resolve texture discriminability, while later-stage VRWKV blocks model long-range dependencies via linear attention to address spatial heterogeneity. This design combines the texture sensitivity of CNNs, the generalization of MLPs, and the scalability of VRWKV. By bridging local and global processing, DualSeg achieves robust multi-scale mapping of morphological priors. As demonstrated in Fig. 1, our model achieves a global, clean ERF (IV) and delivers SOTA performance with optimal computational efficiency (Bottom).

**Page 1, Figure 1:**



**Fig. 1: Top:** Visualization of the *Effective Receptive Fields* (ERF) for different architectures. CNNs (I) focus locally, while Transformers (II) capture global but noisy patterns. Our method (IV) achieves a clean, global ERF. **Bottom:** Performance vs. FLOPs comparison. DualSeg (top-left) achieves the optimal trade-off between segmentation accuracy and computational efficiency.

- 2) **the generalization capability of the model remains insufficiently evaluated.** The experiments are limited to PAS-stained mouse and human datasets that share similar imaging conditions. Without cross-stain or cross-center validation, the robustness of the model under real-world variations in staining or scanner parameters cannot be confirmed. The authors should include additional experiments or at least discuss the expected behavior of the model under stain variability. It is also recommended to cite and discuss two closely related works—“Unsupervised stain augmentation enhanced glomerular instance segmentation on pathology images” and “Identifying and matching 12-level multistained glomeruli via deep learning for diagnosis of glomerular diseases”—to better position this study within the current research landscape.

**Response:** Thank you for your constructive suggestions. We have **incorporated the KPMP dataset** for a more comprehensive evaluation of generalizability and **updated our reference**, as detailed below:

- **Expanded Validation Experiments:** We have incorporated a cross-center validation by performing direct inference on the **KPMP** dataset using models trained exclusively on KPIs data, which is detailed in **Section IV. A. Datasets** and results are illustrated in **Table III and Figure 9**. This evaluation demonstrates DualSeg’s exceptional stability across significant biological and technical variations without the need for domain adaptation.
- **Expanded Discussion on Staining Variability:** While public datasets for non-PAS stains remain scarce, we have added a dedicated discussion on stain variability and its impact on clinical deployment in **Section VI. D. Limitations and Future Work**.
- **Updated Literature and Citations:** we have integrated and discussed the suggested literature in **Section II. B. Technological Evolution of Glomerular Segmentation Architectures**, which better positions DualSeg within the current landscape of robust renal histopathology analysis.

The revisions can be found in Section II. B. Related Work, Section IV. A. Dataset, Section VI. D. Limitations and Future Work, Table III and Figure 9:

**Page 3, Section II. B. Related Work:**

To bridge the gap between local precision and global context, Transformer-based methods like SegFormer [11] introduced self-attention mechanisms, the effectiveness of which has been corroborated in glomerular segmentation studies [37]–[39]. However, standard self-attention faces scalability constraints when processing the extensive spatial dimensions characteristic of Whole Slide Images (WSIs). This limitation has catalyzed the emergence of alternative global modeling paradigms—notably VRWKV [23] and U-Mamba [25], [40]—which utilize recurrent formulations or SSM to achieve effective global receptivity on high-resolution inputs. Although existing hybrid frameworks attempt to synergize the strengths of convolution and attention mechanism [13], achieving seam less multi-scale integration that fully preserves structural continuity remains an ongoing challenge.

**Page 8, Section IV. A. Datasets:**

Dataset III: Human Glomeruli (KPMP). To assess cross species generalization, we retrieved a second human dataset from the Kidney Precision Medicine Project (KPMP) Atlas Repository [50]. Four PAS-stained SVS format WSIs (avg. resolution  $84,000 \times 50,000$ ) were selected with corresponding masks. To rigorously validate generalization, models trained solely on the mouse KPIs dataset were directly applied to this human dataset without retraining. For preprocessing consistency, KPMP WSIs were partitioned into  $2,048 \times 2,048$  patches.

## Page 12, Section VI. D. Limitations and Future Work:

While DualSeg demonstrates superior performance in glomerular segmentation, three primary limitations remain to be addressed in future iterations. First, the model's generalization to rare pathological subtypes, such as global glomerulosclerosis, is currently constrained by data scarcity. We plan to mitigate this by employing domain adaptation and transfer learning techniques to enhance feature robustness for these underrepresented classes. Second, our current validation is restricted to PAS-stained images. To ensure broad clinical applicability, we will extend our evaluation to include Hematoxylin-Eosin (HE) stained datasets, verifying the model's adaptability to varying histological protocols. Finally, despite the linear complexity of the VRWKV block, the computational overhead for gigapixel WSI processing remains significant. Future work will focus on model quantization and lightweight optimization to facilitate deployment on resource-constrained platforms and edge devices.

## Page 7, Table III:

TABLE III: CROSS-DATASET INFERENCE PERFORMANCE COMPARISON FOR GLOMERULAR SEGMENTATION ON THE KPMP DATASET USING 5-FOLD MOUSE-TRAINED MODELS WITH RESPECT TO EXISTING METHODS

Models	1				2				3				4				AVG			
	mDSC $\downarrow$	HD95 $\downarrow$	IoU $\uparrow$	mDSC $\downarrow$	HD95 $\downarrow$	IoU $\uparrow$	mDSC $\downarrow$	HD95 $\downarrow$	IoU $\uparrow$	mDSC $\downarrow$	HD95 $\downarrow$	IoU $\uparrow$	mDSC $\downarrow$	HD95 $\downarrow$	IoU $\uparrow$	mDSC $\downarrow$	HD95 $\downarrow$	IoU $\uparrow$		
U-Net [7]	0.3936 ±0.0045	279.4223 ±2.4891	0.5095 ±0.0045	0.2662 ±0.0042	471.9885 ±2.5645	0.2560 ±0.0042	0.5189 ±1.8505	119.5596 ±0.0045	0.973 ±0.0045	0.5387 ±0.0045	152.3241 ±2.5045	0.9695 ±0.0045	0.4836 ±0.2460	196.2395 ±0.0046	0.4588 ±0.2406	0.3936 ±0.0045	279.4223 ±2.4891	0.5095 ±0.0045	0.2662 ±0.0042	
Attention U-Net [41]	0.7165 ±0.0041	161.1531 ±2.1667	0.6971 ±0.0041	0.4762 ±0.0041	161.1531 ±2.1667	0.6971 ±0.0041	0.8431 ±0.0033	11.5640 ±1.1642	0.973 ±0.0034	0.6220 ±0.0034	11.7108 ±1.1642	0.9650 ±0.0034	0.7056 ±0.0034	12.0369 ±1.1642	0.9736 ±0.0034	0.7160 ±0.0041	161.1531 ±2.1667	0.6971 ±0.0041	0.4762 ±0.0041	
SegNext [26]	0.7926 ±0.0036	164.5795 ±2.2619	0.7056 ±0.0037	0.5235 ±0.0036	86.5292 ±2.3992	0.4978 ±0.0046	0.4848 ±0.0046	36.7357 ±0.0046	0.7287 ±0.0033	0.5685 ±0.0033	30.8195 ±1.7569	0.9048 ±1.7569	0.7410 ±1.7569	112.1102 ±1.7569	0.7116 ±1.7569	0.5235 ±0.0040	86.5292 ±2.3992	0.7056 ±0.0037	0.5235 ±0.0036	
DeepLabV3+ [34]	0.7846 ±0.0034	100.0552 ±2.1076	0.7640 ±0.0034	0.6140 ±0.0034	182.8433 ±2.4553	0.5615 ±0.0047	0.8545 ±0.0047	67.7374 ±1.642	0.8333 ±0.0034	0.6952 ±0.0034	116.7432 ±1.642	0.6617 ±1.642	0.7317 ±1.642	112.0154 ±1.642	0.7018 ±1.642	0.6140 ±0.0040	182.8433 ±2.4553	0.5615 ±0.0047	0.6140 ±0.0034	
Wave-MLP [22]	0.8152 ±0.0034	145.4421 ±2.7809	0.7580 ±0.0036	0.4555 ±0.0036	4.0229 ±2.7413	0.4029 ±0.0047	0.8505 ±0.0047	60.6731 ±0.8149	0.8253 ±0.0033	0.7675 ±0.0033	92.2304 ±1.6319	0.7342 ±1.6319	0.7046 ±1.6319	106.0041 ±1.6319	0.7353 ±1.6319	0.4555 ±0.0038	4.0229 ±2.7413	0.4029 ±0.0047	0.8505 ±0.0046	
InceptionNet [42]	0.6564 ±0.0044	236.3450 ±2.3635	0.6352 ±0.0045	0.4779 ±0.0045	272.6440 ±2.6615	0.4672 ±0.0045	0.7136 ±0.0045	167.7567 ±0.9100	0.6937 ±0.0045	0.4303 ±0.0045	312.1570 ±2.4703	0.8875 ±2.4703	0.5279 ±2.4703	284.8018 ±2.4428	0.4967 ±0.0045	0.6352 ±0.0046	272.6440 ±2.6615	0.4672 ±0.0045	0.4779 ±0.0045	
SegFormer [11]	0.8037 ±0.0034	125.7612 ±2.4737	0.7814 ±0.0034	0.6451 ±0.0034	206.3655 ±2.6615	0.8515 ±0.0045	0.5897 ±0.0045	47.1347 ±0.9100	0.8353 ±0.0045	0.8285 ±0.0045	49.0933 ±1.9100	0.7986 ±1.9100	0.8082 ±1.9100	70.7543 ±1.9100	0.7802 ±1.9100	0.6451 ±0.0045	206.3655 ±2.6615	0.8515 ±0.0045	0.5897 ±0.0045	
VRWKV [23]	0.7912 ±0.0036	143.4809 ±2.1290	0.7662 ±0.0037	0.5112 ±0.0036	271.6789 ±2.2767	0.4883 ±0.0047	0.8056 ±0.0047	62.1875 ±1.0115	0.7650 ±0.0037	0.6937 ±0.0037	76.0115 ±1.3855	0.7296 ±1.3855	0.7480 ±1.3855	80.0189 ±1.3855	0.7188 ±1.3855	0.5112 ±0.0039	271.6789 ±2.2767	0.4883 ±0.0047	0.8056 ±0.0046	
VM-UNet-V2 [40]	0.6521 ±0.0043	181.1119 ±2.1959	0.6233 ±0.0044	0.4833 ±0.0044	332.8376 ±2.5972	0.4575 ±0.0046	0.7880 ±0.0046	104.5297 ±1.5297	0.7620 ±0.0046	0.5410 ±0.0046	232.7710 ±2.0043	0.6063 ±2.0043	0.6028 ±2.0043	216.0202 ±2.3843	0.5664 ±2.3843	0.4833 ±0.0044	332.8376 ±2.5972	0.4575 ±0.0046	0.7880 ±0.0045	
UNETR [9]	0.5946 ±0.0046	99.0704 ±2.1782	0.5670 ±0.0047	0.4500 ±0.0047	301.7408 ±2.5052	0.4410 ±0.0047	0.8059 ±0.0047	154.2609 ±1.8669	0.8059 ±0.0047	0.8059 ±0.0047	333.3864 ±2.0834	0.6499 ±2.0834	0.6049 ±2.0834	303.0103 ±2.2676	0.5947 ±2.2676	0.4500 ±0.0047	301.7408 ±2.5052	0.4410 ±0.0047	0.8059 ±0.0046	
Swin UNETR [43]	0.5171 ±0.0047	207.6704 ±2.4881	0.4921 ±0.0046	0.3438 ±0.0046	432.7990 ±2.6615	0.3300 ±0.0045	0.7037 ±0.0045	85.6554 ±0.9100	0.6805 ±0.0045	0.6559 ±0.0045	131.7022 ±1.9244	0.6201 ±1.9244	0.6137 ±1.9244	156.2208 ±2.2260	0.5824 ±2.2260	0.4921 ±0.0044	432.7990 ±2.6615	0.3300 ±0.0045	0.7037 ±0.0046	
DA-TransM-Net [44]	0.7783 ±0.0039	97.7011 ±2.4737	0.7563 ±0.0039	0.5907 ±0.0039	166.4438 ±2.6615	0.5610 ±0.0045	0.8319 ±0.0045	63.6652 ±0.9100	0.8143 ±0.0045	0.7948 ±0.0045	81.5234 ±1.9100	0.7594 ±1.9100	0.7780 ±1.9100	87.6505 ±1.9100	0.7489 ±1.9100	0.5907 ±0.0045	166.4438 ±2.6615	0.5610 ±0.0045	0.8319 ±0.0045	
H2Former [10]	0.7490 ±0.0039	146.2938 ±2.7220	0.7220 ±0.0040	0.5627 ±0.0040	56.0562 ±2.3340	0.5005 ±0.0045	0.8095 ±0.0045	55.2729 ±1.7559	0.5639 ±0.0045	0.6939 ±0.0045	169.0482 ±2.0042	0.5951 ±2.0042	0.6815 ±2.0042	242.2960 ±2.3043	0.6460 ±2.3043	0.5627 ±0.0042	56.0562 ±2.3340	0.5005 ±0.0046	0.8095 ±0.0045	
U-mamba [25]	0.6549 ±0.0043	132.5547 ±2.1727	0.6214 ±0.0043	0.3967 ±0.0043	315.6740 ±2.7300	0.3729 ±0.0040	0.7500 ±0.0040	75.1434 ±1.5002	0.7262 ±1.5002	0.5380 ±1.5002	192.8970 ±2.0043	0.5878 ±2.0043	0.5843 ±2.0043	181.0034 ±2.1727	0.5449 ±2.1727	0.3967 ±0.0043	315.6740 ±2.7300	0.3729 ±0.0040	0.7500 ±0.0041	
DualSeg(ours)	<b>0.8251</b> ±0.0034	<b>81.6201</b> ±1.5314	<b>0.8022</b> ±0.0035	<b>0.5888</b> ±0.0045	<b>235.8327</b> ±2.6482	<b>0.5573</b> ±0.0045	<b>0.8848</b> ±0.0028	<b>21.0608</b> ±0.2919	<b>0.8630</b> ±0.0029	<b>0.8393</b> ±0.0029	<b>57.6049</b> ±1.3786	<b>0.8123</b> ±1.3786	<b>0.8195</b> ±1.3786	<b>69.6369</b> ±1.5420	<b>0.7938</b> ±1.5420	<b>0.8251</b> ±0.0034	<b>81.6201</b> ±1.5314	<b>0.8022</b> ±0.0035	<b>0.5888</b> ±0.0045	

## Page 10, Figure 8:

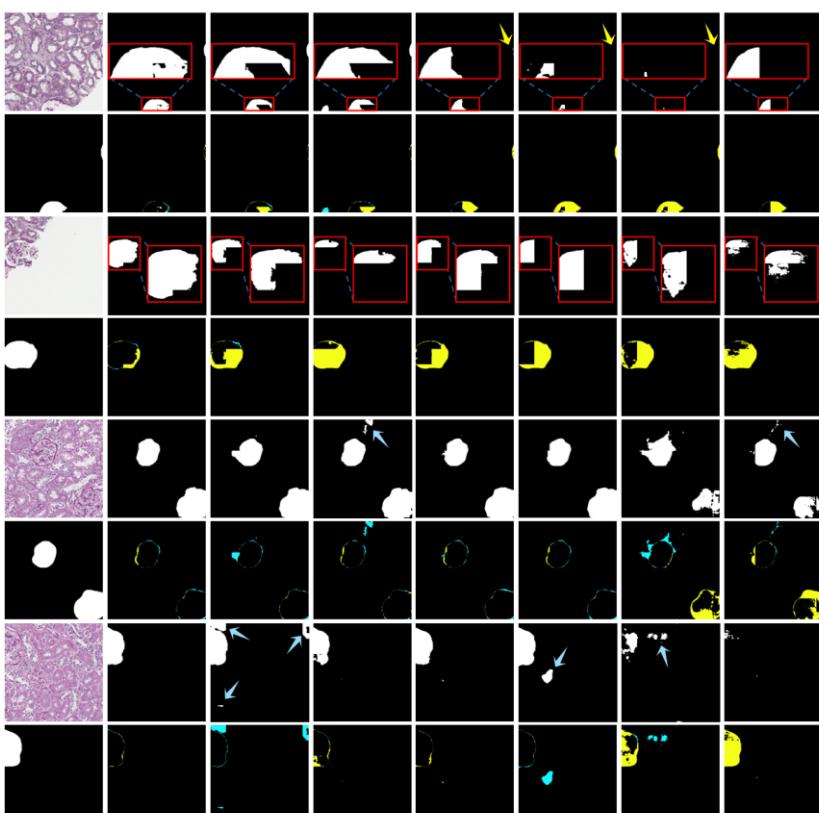


Fig. 8: Visual comparison on the held-out KPMP test set. Odd rows display original images and inference masks; even rows show GT and error maps, where yellow and green denote under-segmentation and over-segmentation, respectively. Red boxes and matching arrows highlight magnified details and segmentation errors. DualSeg exhibits unprecedented stability when handling cross-center and cross-species data.

- 3) the clinical relevance and interpretability of the segmentation results are not adequately discussed. Although the paper emphasizes segmentation accuracy, it does not explain how this improvement translates into clinical or diagnostic benefits, such as better lesion quantification, assessment of glomerulosclerosis, or prediction of disease progression. Including examples or quantitative analyses connecting segmentation quality to potential diagnostic metrics would significantly strengthen the manuscript. A brief discussion of the model's interpretability from a pathologist's perspective would also be valuable.

**Response:** Thank you for the feedback. We have revised the manuscript to bridge the gap between technical accuracy and diagnostic utility, as detailed below:

- **Statistical Validation of Clinical Utility:** We have enriched **Section VI. C. Clinical Relevance** with statistical analysis to demonstrate the diagnostic value of our mode. Specifically, we incorporated a paired *t*-test analysis in **Figure 11**, demonstrating that DualSeg's performance gains are statistically significant ( $p < 0.05$  to  $p < 0.001$ ) across diverse cohorts. This statistical robustness supports the model's reliability for the standardized quantification of critical biomarkers, such as glomerulosclerosis and fibrosis, which are essential for predicting disease progression.
- **Adding Error Map:** To enhance interpretability from a pathologist's perspective, we have introduced Error Maps in **Figures 7 and 9**. These maps utilize red bounding boxes to magnify local anatomical details and provide a spatial audit of model performance by color-coding under-segmentation (yellow) and over-segmentation (green). This visualization allows clinicians to assess diagnostic confidence more effectively and helps mitigate inter-observer variability in complex histopathological scenarios .

The revisions can be found in Section VI. C. Clinical Relevance, Figures 7, 9 and 11:

**Page 13, Section VI. C. Clinical Relevance:**

DualSeg exhibits statistically significant superiority ( $p < 0.05$ – $0.001$ ; Fig. 11) and exceptional reproducibility, evidenced by a minimal standard deviation (0.0040) on the HuBMAP dataset. Its ability to accurately resolve diverse morphologies—ranging from mild hypertrophy to severe fragmentation—enables the precise quantification of pathological biomarkers like sclerosis and fibrosis. Furthermore, the model's robust performance on the cross-species KPMP dataset supports standardized CKD monitoring. By mitigating inter-observer variability and reducing manual annotation burdens, DualSeg provides a scalable solution for multi-center clinical trials and routine diagnostic workflows.

Page 8, Figure 7:

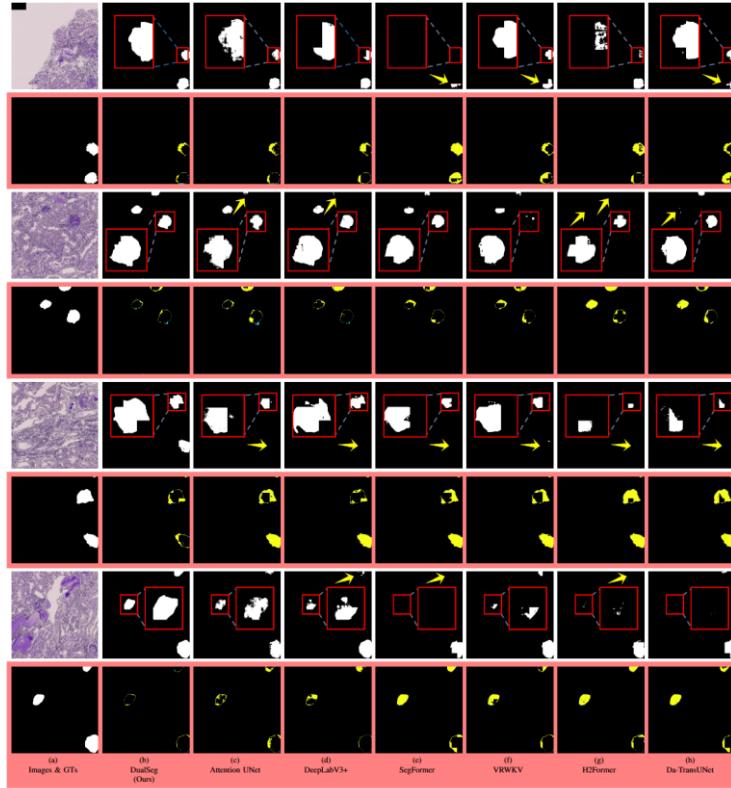


Fig. 7: Visual comparison on the KPIs test set. **Odd rows** display original images and inference masks; **even rows** show GT and error maps, where **yellow** and **green** denote under-segmentation and over-segmentation, respectively. Red boxes and matching arrows highlight magnified details and specific segmentation errors. **DualSeg** demonstrates superior performance in handling heterogeneous glomeruli with intricate boundaries.

Page 11, Figure 9:

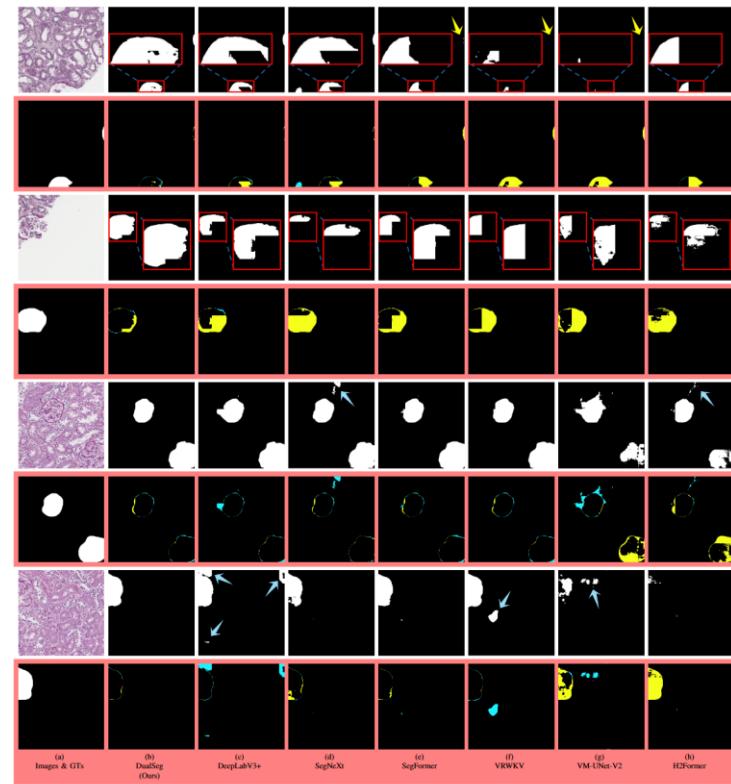


Fig. 9: Visual comparison on the held-out KPMP test set. **Odd rows** display original images and inference masks; **even rows** show GT and error maps, where **yellow** and **green** denote under-segmentation and over-segmentation, respectively. Red boxes and matching arrows highlight magnified details and segmentation errors. **DualSeg** exhibits unprecedented stability when handling cross-center and cross-species data.

Page 12, Figure 11:

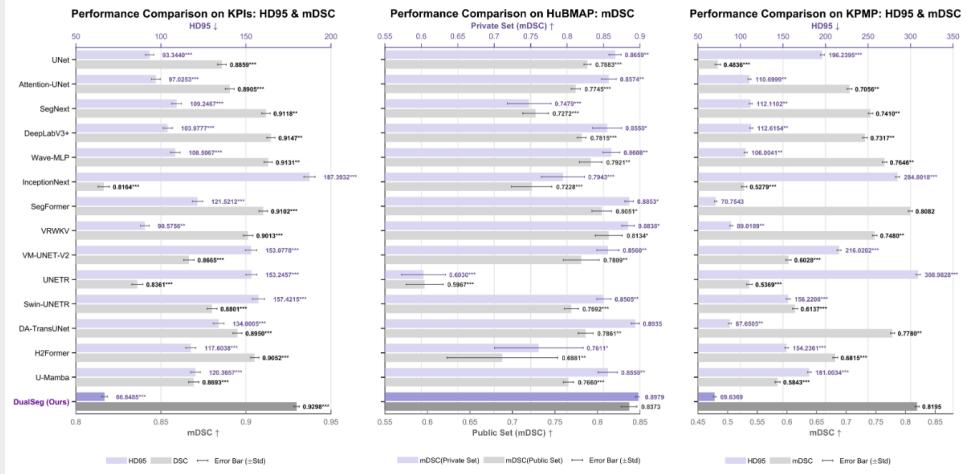


Fig. 11: Performance comparison between our DualSeg model and 14 baseline methods across three datasets. The comparison metrics include mDSC and HD95 for the KPIs and KPMP datasets (left and right panels, respectively), and mDSC for the HuBMAP dataset (middle panel). The error bars represent  $\pm$  standard deviation. Statistical significance was assessed using paired *t*-tests, with levels indicated by asterisks: \* $p < 0.05$ , \*\* $p < 0.01$ , \*\*\* $p < 0.001$ .

- 4) **the methodological presentation is unnecessarily complicated.** The mathematical description of the Wave-Swin and VRWKV blocks is long and dense, making it difficult for readers from the biomedical community to grasp the main idea. Some mathematical symbols, such as  $\Theta$  and  $W^T_{jk}$ , are not clearly defined, and equation formatting is occasionally inconsistent. The authors should simplify the explanation of equations and focus on the conceptual understanding of each component's role in the overall framework, leaving detailed derivations to supplementary materials if necessary.

**Response:** Thank you for your advice. We have streamlined the methodological presentation in **Section III** to improve accessibility for the biomedical community while maintaining technical precision. The mathematical descriptions of the Wave-Swin and VRWKV blocks have been simplified to focus on their conceptual roles in resolving texture discriminability and spatial heterogeneity. We have standardized the formatting across all mathematical expressions—specifically **Equations (1) through (5)**—and provided explicit definitions for all symbols to ensure clarity and consistency. By prioritizing the structural intuition of the dual-stage framework, we have ensured the methodology remains both rigorous and accessible to a broader audience.

The revisions can be found in Section III Methodology:

#### Page 4, Section III Methodology:

##### A. Wave-Swin Block

...

1) *Wave Formulation:* Let the input feature map be denoted as  $X = [x_1, x_2, \dots, x_n]$ , where each  $x_j$  represents a token. Instead of standard linear projections, we map these tokens into a complex wave representation to capture both semantic intensity and spatial structural priors. We define the complex wave form  $z_j$  for the  $j$ -th token as:

$$z_j = \mathcal{A}(x_j) \odot \exp(i \cdot \mathcal{P}(x_j)), \quad (1)$$

where  $\mathcal{A}(\cdot)$  and  $\mathcal{P}(\cdot)$  denote learnable linear transformations that project the input to amplitude and phase terms, respectively. The operator  $\odot$  represents element-wise multiplication. ... These components are then aggregated via a dynamic token mixing operation:

$$y_j = \sum_{k \in \Omega_j} (\text{Re}(z_k) W_{Re} + \text{Im}(z_k) W_{Im}), \quad (2)$$

Where  $y_j$  is the output token,  $\Omega_j$  denotes the dynamic propagation window centered at  $j$ , and  $W_{Re}$ ,  $W_{Im}$  are learnable weights governing the fusion of spatial components.

...

##### B. VRWKV Block

...

1) *Z-Shift and Spatial Mixing:* ...

After the spatial shift, the feature maps are flattened into token sequences and projected. The generation of the receptor, key, and value matrices is formulated as:

$$N_s = \text{Linear}_N (\operatorname{Flatten}(\text{Z-Shift}(X))), \quad N \in \{R, K, V\}, \quad (3)$$

where  $X$  denotes the input 2D feature map, and  $\text{Linear}_N$  represents the learnable projection weights. The flattened tokens then undergo the linear-complexity bidirectional attention aggregation:

$$S = \sigma(R_s) \odot \text{Bi-WKV}(K_s, V_s), \quad (4)$$

where  $\sigma(\cdot)$  is the sigmoid activation,  $\odot$  denotes element-wise multiplication, and Bi-WKV is the time-mixing operator that aggregates global context with linear complexity  $O(L)$ , efficiently capturing pairwise affinities between distant glomerular candidates.

2) *Channel Mixing*: Following spatial aggregation, the features undergo Channel Mixing to enable inter-channel communication. This module mirrors the gating mechanism of the spatial stage but focuses on feature refinement within the channel dimension. The transition is expressed as:

$$O_c = \sigma(R_c) \odot (\text{SqReLU}(K_c) \cdot W_v), \quad (5)$$

where  $R_c$  and  $K_c$  are derived from the spatially mixed features via linear projections, and SqReLU denotes the squared ReLU activation.

**5) the experimental analysis requires stronger statistical and methodological support.**

Although the authors conduct ablation studies, the reported improvements are small, and the absence of variance analysis or statistical testing makes it unclear whether the gains are significant. The paper would benefit from a more comprehensive evaluation, including inference time, parameter count, and performance on challenging subtypes such as sclerotic or crescentic glomeruli. Additional qualitative examples demonstrating both the strengths and failure cases of DualSeg would help provide a balanced assessment of the model's robustness.

**Response:** Thank you for your advice. We have included a summary of t-tests, FLOPs comparison, and failure case analysis to strengthen the empirical foundation of our study. We have addressed the specific concerns as follows:

- **Enhanced Statistical Analysis:** We have significantly strengthened the statistical and methodological support for our experimental analysis. As detailed in newly added **Section VI. C. Clinical Relevance** and **Figure 11**, we performed **t-tests across all three datasets**, confirming that DualSeg's performance gains are statistically significant ( $p < 0.05$  to  $p < 0.001$ ). To address reproducibility, we included variance analysis (error bars) in all primary results, with DualSeg exhibiting exceptional stability (e.g., a minimal standard deviation of 0.0040 on the HuBMAP dataset).
- **Expanded Efficiency Evaluation:** We incorporated a **performance-vs-FLOPs comparison in Figure 1 and Figure 3** to justify computational efficiency and expanded our qualitative evaluation to include challenging subtypes such as fragmented and sclerotic glomeruli.
- **Addition of Failure Case Assessment:** Finally, we have provided a balanced assessment in **Section VI. B. Failure Case Analysis and Figure 10** by analyzing specific failure cases, offering a transparent view of the model's current limitations and future improvement directions.

The supporting data for these revisions can be found in Figure 1, 3(pertaining to Comment 1), 11, Section VI. C. Clinical Relevance (pertaining to Comment 3), Figure 10 and Section VI. B. Failure Case Analysis, which are provided below:

**Page 12, Section VI. B. Failure Case Analysis:**

Fig. 11 reveals that the model occasionally fails to detect globally sclerotic glomeruli in the KPMP dataset. This limitation stems primarily from two factors: the partial truncation of peripheral glomeruli during WSI tiling, which compromises morphological context, and the significant divergence of unseen, extreme pathological variants. To mitigate this, future work could increase patch sizes to preserve boundary information or, more efficiently, integrate uncertainty-guided semi-supervised learning. This strategy aims to enhance robustness against rare phenotypes without incurring excessive computational overhead.

**Page 12, Figure 10:**

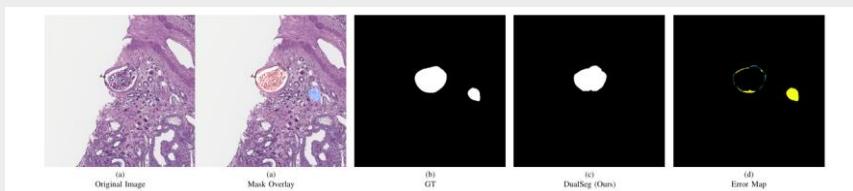


Fig. 10: Visual analysis of segmentation limitations. In the error map (far right), the yellow region highlights a significant false negative (under-segmentation). This failure is primarily attributed to the absence of globally sclerotic samples in the training set, preventing the model from generalizing to the *high heterogeneity and distinct morphological features* of this unseen pathology.

**Finally, several minor issues should be corrected:**

- 1) Figure fonts are too small to read, and color schemes in Figures 3–6 are inconsistent. Figure 5 panels are misaligned, and the arrows indicating regions of interest are too faint.
- 2) Table I has several misaligned columns and overlapping text.
- 3) There are typographical errors such as “uqualitative” (should be “qualitative”) and repeated line-break hyphenations such as “glomeru- lar,” which should be corrected.
- 4) Ensure consistent capitalization of dataset names (e.g., HuBMAP, KPIs) and verb tense consistency throughout the text. References should be carefully checked for completeness and formatting.

**Response:** We have performed a comprehensive editorial overhaul of the manuscript:

- **Figure Enhancements:** We have revised all figures to ensure maximum legibility and stylistic consistency. Specifically, font sizes were increased across all plots, and the color schemes in **Figures 3–6** were unified to maintain a cohesive visual identity. Panel alignment in **Figure 5** has been corrected, and the indicating arrows have been thickened and brightened to clearly highlight regions of interest.

The revisions can be found in Figure 5 and 6:

**Page 4, Figure 5:**

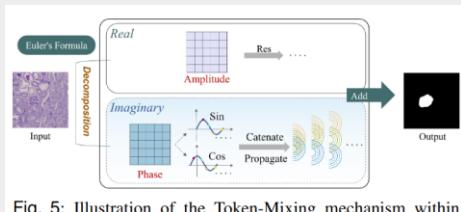


Fig. 5: Illustration of the Token-Mixing mechanism within the **Wave-Swin Block**. The input is initially decomposed into a *real* component and an *imaginary* component via Euler's formula. The modulated phase undergoes two distinct cosine transformations and window-based propagation, and is subsequently fused with the amplitude information to generate the output.

**Page 5, Figure 6:**

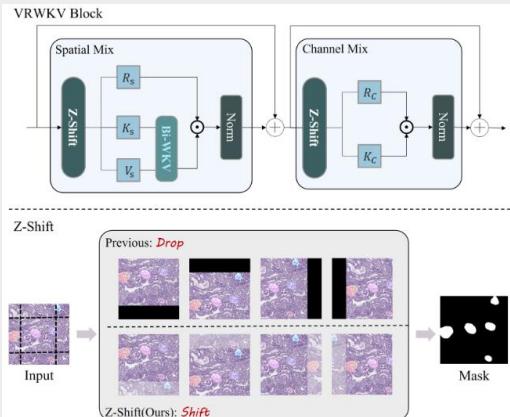


Fig. 6: Structure of the improved VRWKV block. The black margins in the lower panel illustrate the prior approach where pixels shifted beyond the boundary are simply dropped. In contrast, our **Z-Shift** wraps these pixels to the opposite edge, filling the empty regions and eliminating information loss.

- **Table Reformatting:** Table I has been reformatted to resolve column misalignments and overlapping text. The presentation of metrics, including mDSC, HD95, and IoU, is now clear and professionally aligned .
- **Textual Corrections:** All typographical errors, such as "uqualitative," have been corrected. We have also removed improper line-break hyphenations (e.g., "glomeru-lar") throughout the text to ensure linguistic fluidity.
- **Consistency and References:** We have conducted a full audit of the manuscript to ensure consistent capitalization of dataset names, such as **HuBMAP, KPIs, and KPMP**. Verb tenses have been unified, and the reference list has been verified for completeness and adherence to the journal's formatting standards.