

Bridging RGB and Hematoxylin Components: An Interleaved Guidance and Fusion Framework for Point Supervised Nuclei Segmentation

Supplementary Material

001 1. Discussion and Visualization on IGA

002 To further examine the contribution of the proposed Inter-
003 leaved Guidance Attention (IGA), we visualize four compo-
004 nents: the RGB input, the ground-truth instance mask (GT),
005 the raw prediction maps with IGA (IGA), and those with-
006 out IGA (NIGA). As shown in Figure 1, incorporating IGA
007 substantially enhances the model’s ability to focus on nuclei
008 regions under weak supervision.

009 Compared with the RGB input, the IGA responses sup-
010 press irrelevant background textures and staining noise, en-
011 abling the network to concentrate on coherent morpholog-
012 ical structures. When aligned with GT, the IGA predic-
013 tion maps exhibit stronger spatial compactness and clearer
014 boundaries, reflecting a more faithful reconstruction of nu-
015 clear shapes and a better separation between adjacent in-
016 stances. In contrast, NIGA predictions often suffer from
017 dispersed activations, spurious responses on background re-
018 gions, and incomplete emphasis on nuclei with irregular
019 morphology.

020 These observations highlight the essential role of IGA
021 in establishing robust cross-representation guidance. By
022 enforcing complementary semantic cues and spatial con-
023 sistency, IGA encourages the network to refine activations
024 at nuclei centers while imposing additional structural con-
025 straints along boundaries. This mechanism effectively miti-
026 gates the instability caused by sparse point supervision and
027 label noise, yielding more focused, consistent, and separa-
028 ble nuclei representations. Overall, IGA significantly im-
029 proves the reliability of subsequent instance segmentation.

030 2. Efficiency and Performance Analysis

031 As shown in Figure 2, our method (Final) achieves the
032 best overall performance across multiple metrics while
033 maintaining moderate computational complexity. On the
034 MO dataset [2], it records the highest Accuracy (91.38%),
035 F1-score (79.03%), Dice_{obj} (74.89%), and AJI (54.85%),
036 demonstrating the effectiveness of our dual-representation
037 fusion mechanism in enhancing instance-level nuclei seg-
038 mentation quality.

039 Compared with the lightweight D2E2 [7], our method
040 yields substantial improvements, including a +3.9% gain in
041 F1-score and a +4.9% increase in AJI. Although it requires
042 considerably fewer FLOPs and parameters than the more
043 complex SCNet [3] (46.18 G rather than 73.50 G FLOPs,
044 and 65.1 M rather than 80.1 M parameters), it still surpasses
045 SCNet across all evaluation metrics, representing a 37% re-

Train→Test	Method	Acc	F1-score	Dice _{obj}	AJI
M→S	Ours	90.23 _{0.21}	65.59 _{0.18}	60.97 _{0.22}	36.30 _{0.19}
	MA	88.95 _{0.37}	62.42 _{0.41}	58.15 _{0.33}	34.21 _{0.48}
S→M	Ours	88.57 _{0.25}	67.62 _{0.20}	56.34 _{0.28}	38.35 _{0.24}
	MA	86.70 _{0.46}	64.75 _{0.39}	51.91 _{0.44}	32.80 _{0.40}
M→C	Ours	87.96 _{0.17}	69.83 _{0.23}	67.62 _{0.19}	48.24 _{0.22}
	MA	89.45 _{0.35}	71.02 _{0.44}	69.85 _{0.38}	49.13 _{0.41}
C→M	Ours	88.65 _{0.26}	69.98 _{0.22}	66.89 _{0.25}	46.37 _{0.20}
	MA	84.36 _{0.49}	66.15 _{0.34}	64.78 _{0.42}	43.92 _{0.37}
C→S	Ours	86.28 _{0.19}	48.87 _{0.27}	47.82 _{0.24}	25.86 _{0.23}
	MA	83.96 _{0.45}	40.65 _{0.38}	44.76 _{0.41}	22.51 _{0.36}
S→C	Ours	86.38 _{0.23}	63.15 _{0.21}	60.88 _{0.28}	39.25 _{0.26}
	MA	85.18 _{0.33}	60.21 _{0.40}	54.50 _{0.43}	35.50 _{0.39}

Table 1. Adaptive domain generalization results across three datasets. Each training-testing pair is evaluated using two methods: our adaptive fusion approach (Ours) and MixAnno (MA). Evaluation metrics include Accuracy (Acc), F1-score (F1), Dice_{obj}, and Aggregated Jaccard Index (AJI). M, S, and C denote the MO, CoNSeP, and CPM-17 datasets, respectively.

duction in computational cost while delivering superior ac- 046
curacy. 047

Relative to the mixed-annotation method MA [5], our 048
approach achieves slightly higher accuracy (91.38% com- 049
pared to 91.21%) and produces notable improvements in 050
F1-score (+2.1%) and AJI (+3.2%), reflecting stronger ro- 051
bustness and generalization under point level supervision. 052

Overall, the proposed framework provides a well- 053
balanced trade-off between segmentation precision and 054
computational efficiency, which can be attributed to its 055
cross-representation complementarity and adaptive fusion 056
strategy. 057

058 3. Domain Generalization Performance

Table 1 presents the cross-domain generalization results 059
among MO (M) [2], CoNSeP (S) [1], and CPM-17 060
(C) [6]. Our adaptive fusion framework consistently 061
demonstrates superior or comparable performance across 062
all training-testing configurations, highlighting its robust- 063
ness to domain shifts in histopathological images. Specif- 064
ically, when transferring from MO to CoNSeP (M→S), 065
our method surpasses MixAnno [5] by +2.8% in F1-score 066
and +2.8% in Dice, reflecting improved structural align- 067
ment under staining and morphological variations. Simi- 068
larly, under the reverse setting (S→M), we observe gains 069
of +2.9 F1% and +4.4 Dice%, confirming enhanced cross- 070

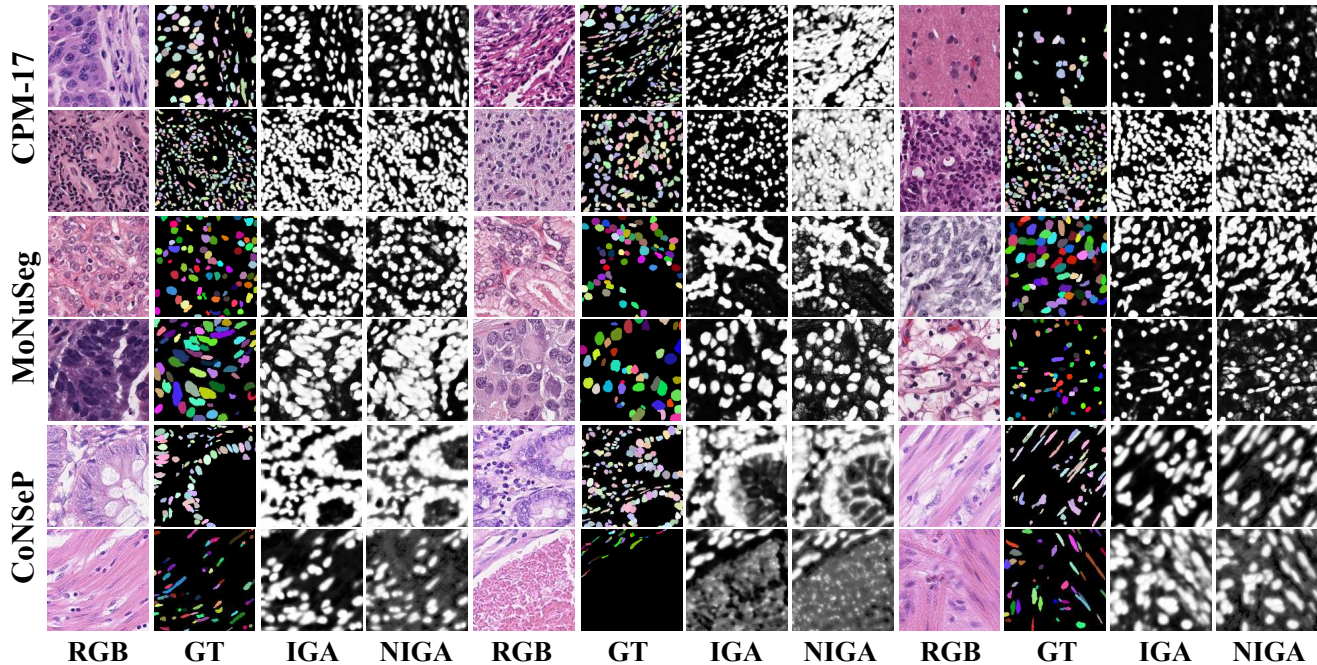


Figure 1. Visualization of IGA across the MoNuSeg (MO), CPM-17, and CoNSEP datasets. For each dataset, the four columns show: (1) RGB input image; (2) ground-truth (GT) instance annotations; (3) IGA prediction logits generated with the proposed IGA module; and (4) NIGA prediction logits obtained without IGA.

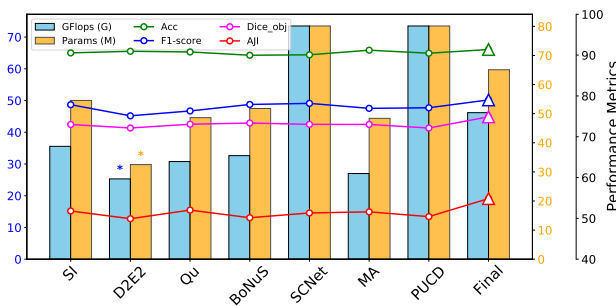


Figure 2. Comparison of different methods in terms of segmentation metrics, parameter counts, and computational complexity on the MO dataset. The left y -axis represents **GFLOPs** (blue), while the right y -axis indicates the number of **parameters** (orange). The best metric results are marked with a **triangle**, and the most efficient methods in terms of GFLOPs and parameters are highlighted with a **star**.

071 domain representation learning. Even in the most challeng-
 072 ing transfer ($C \rightarrow S$), our approach outperforms MixAnno
 073 by a significant margin (**+8.2% in F1-score** and **+3.1%**
 074 **in Dice_{obj}**), indicating effective adaptation to unseen data
 075 distributions. Although MixAnno slightly excels in $M \rightarrow C$
 076 due to dataset-scale bias, our model maintains competi-
 077 tive results and exhibits more stable generalization, as evi-
 078 denced by the smaller standard deviations. Overall, these

results demonstrate that the proposed interleaved adaptive
 fusion strategy effectively leverages complementary cues
 between RGB and Hematoxylin representations, enabling
 strong domain generalization across diverse pathological
 datasets without retraining.

4. Complementarity Analysis

As discussed in the introduction, there are four key issues
 we will analyze: excessive staining, staining artifacts, "Two
 or one nuclei? ", and the challenge of separating nuclei
 from the background. These issues highlight the comple-
 mentarity between RGB and H-channel methods, as each
 method has its own strengths and limitations. In this sec-
 tion, we analyze these issues in detail through the use of
 RGB and H-channel images.

4.1. Excessive Staining

Figure 3(B) illustrates the issue of excessive staining. In
 the original RGB images, over-staining leads to blurred cell
 boundaries or background regions being dyed, which causes
 incorrect segmentation of nuclei. Grayscale images, deriv-
 ed from RGB images, do not significantly improve upon
 this problem as they still suffer from over-staining effects,
 making it difficult to differentiate between nuclei and back-
 ground. In contrast, the H-channel extraction method fil-
 ters out background regions and enhances the contrast be-
 tween the nuclei and background, improving nuclei separa-

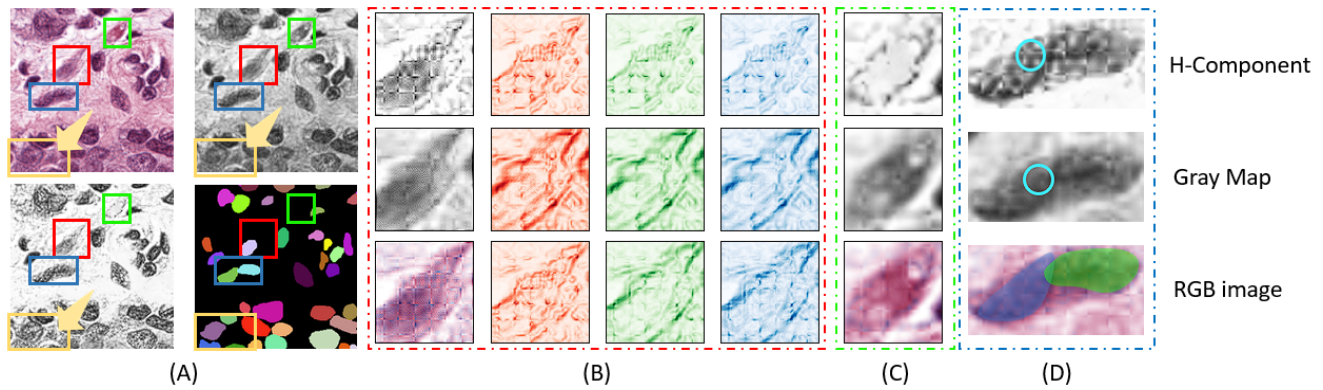


Figure 3. Analysis of image segmentation challenges and H-channel complementarity. (A) Visualizations of different image issues: top-left shows the original RGB image, top-right shows the corresponding grayscale image, bottom-left shows the H-channel image, and bottom-right shows the ground truth labels. The complementarity between RGB and H-channel is highlighted in the yellow box and arrow, particularly in regions where H-channel extraction struggles with background separation in dense nuclear areas. (B) Gradient analysis of excessive staining in RGB images. (C) Staining artifact issue, demonstrating residual staining in non-nuclear regions. (D) The "Two or One Nuclei?" issue, where weak boundaries or insufficient staining lead to difficulties in differentiating overlapping nuclei.

104 tion. However, in areas with severe staining artifacts, the
105 H-channel may still incorrectly classify stained regions as
106 nuclei, leading to segmentation errors.

107 4.2. Staining Artifact

108 Figure 3(C) shows the staining artifact problem. This issue
109 originates from the uneven staining in the original RGB im-
110 age, where non-nuclear areas still exhibit artificial staining.
111 This residual staining impacts nuclei segmentation, mis-
112 classifying background regions as part of the nuclei. The H-
113 channel method, by suppressing areas dominated by Eosin,
114 greatly reduces such artifacts, correctly identifying them as
115 background. Thus, H-channel extraction helps significantly
116 in reducing staining artifacts and improving nuclei segmen-
117 tation accuracy.

118 4.3. Two or One Nuclei?

119 In Figure 3(D), we show the "Two or one nuclei?" prob-
120 lem. In both RGB and grayscale images, weakly stained
121 cell boundaries or insufficient staining may cause nuclei to
122 merge or be hard to distinguish. This results in misclas-
123 sification, such as a single nucleus being misidentified as
124 two or vice versa. Through gradient analysis, we show that
125 the H-channel method enhances the contrast between nuclei
126 and background, assisting in separating overlapping nuclei.
127 However, in cases of severe overlap or unclear boundaries,
128 the H-channel method may still misclassify multiple nuclei
129 as a single instance, leading to blurred nuclei boundaries
130 and inaccurate segmentation.

131 4.4. Background Separation in Dense Nuclei Re- 132 gions

133 Finally, as shown in the yellow box and arrow in Figure
134 3(A), the H-channel method struggles to separate nuclei
135 from the background in densely packed nuclear regions.
136 The staining in these areas may spill over into the back-
137 ground, leading to staining accumulation. In the H-channel
138 image, this accumulation appears as falsely segmented nu-
139 clei, where surrounding regions form what seems like a
140 mound of staining. Therefore, to accurately identify such
141 areas, it is necessary to analyze the original RGB image
142 alongside the H-channel image. This highlights the comple-
143 mentarity between RGB and H-channel methods, as each
144 method excels in different areas and compensates for the
145 other's shortcomings.

146 4.5. Conclusion

147 The RGB and H-channel methods each have their own
148 strengths and weaknesses. RGB images provide rich color
149 information and retain structural context, but are prone to
150 over-staining and background interference, which can cause
151 boundary blurring and staining artifacts. H-channel ex-
152 traction excels at enhancing contrast and suppressing back-
153 ground interference, but it struggles in densely stained re-
154 gions and may misinterpret overlapping nuclei as a single
155 instance. Therefore, the complementarity between RGB
156 and H-channel methods is critical for improving nuclei seg-
157 mentation accuracy. By combining both methods, we can
158 leverage their respective strengths and mitigate their weak-
159 nesses, leading to more robust and accurate segmentation
160 results.

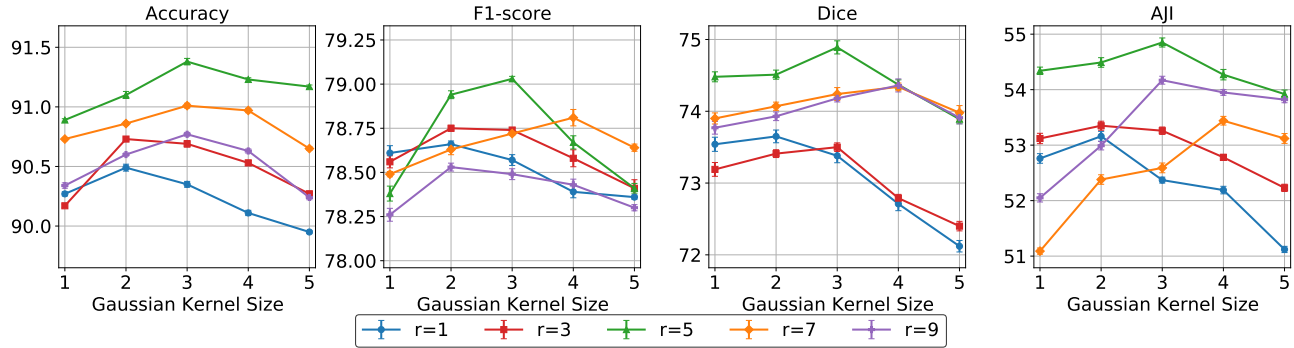


Figure 4. Line chart illustrating the effect of Gaussian kernel radius on four key metrics. Here, r denotes the radius setting, and different colors indicate metric performance under different radius settings.

M	L	S	R	Dice _{obj}	AJI	DQ	SQ	PQ
✓				74.29 _{0.22}	54.41 _{0.19}	69.31 _{0.09}	71.39 _{0.17}	51.19 _{0.24}
	✓			74.38 _{0.18}	54.36 _{0.21}	69.43 _{0.20}	71.25 _{0.15}	51.18 _{0.19}
		✓		74.56 _{0.25}	54.56 _{0.13}	69.42 _{0.17}	71.72 _{0.18}	51.21 _{0.25}
			✓	74.55 _{0.20}	54.58 _{0.16}	69.48 _{0.22}	71.94 _{0.19}	51.46 _{0.26}
✓			✓	74.69 _{0.19}	54.57 _{0.24}	70.33 _{0.15}	72.26 _{0.37}	51.39 _{0.27}
	✓	✓	✓	74.72 _{0.13}	54.69 _{0.10}	70.34 _{0.16}	72.38 _{0.18}	51.38 _{0.25}
✓	✓	✓	✓	74.89 _{0.11}	54.85 _{0.34}	70.65 _{0.17}	72.54 _{0.29}	51.45 _{0.13}

Table 2. Ablation study of each component on the MO dataset. M, L, S, and R denote different modules. Metrics include Dice_{obj}, AJI, DQ, SQ, and PQ (%).

5. Impact of Kernel Size and Radius

As shown in Figure 4, the *Gaussian kernel size* denotes the spatial extent of the kernel applied during Gaussian blurring, while the *point radius* (denoted as r in the legend) specifies the radius assigned to each annotated point prior to blurring.

Empirically, we observe that moderate kernel sizes yield more stable and accurate instance segmentation results. Smaller kernels are insufficient to suppress annotation noise, whereas excessively large kernels tend to oversmooth structural boundaries and degrade fine-grained localization. A kernel size within the range of 3–5 offers an effective balance between noise attenuation and structural fidelity.

Regarding the point radius, increasing r from 1 to 5 consistently enhances segmentation performance by providing richer contextual cues around annotated centers. However, an overly large radius (e.g., $r = 9$) introduces spatial ambiguity, and the localization accuracy slightly deteriorates.

Overall, combining a moderate Gaussian kernel with a relatively large point radius achieves an optimal trade-off between contextual aggregation and spatial precision in weakly supervised nuclei instance segmentation.

6. Ablation Analysis of RCDF Components

Table 2 presents the ablation results of different components within the RCDF module, including the multi-scale feature-aware units (M), the lightweight global context-aware scale modulator (L), the sequential attention (S), and the residual guided fusion (R). Individually adding each component leads to consistent improvements across all metrics, demonstrating their effectiveness in enhancing feature representation and fusion. In particular, introducing the residual guided fusion (R) yields a notable gain in Dice_{obj} (from 74.29 to 74.55) and AJI (from 54.41 to 54.58), indicating that the residual pathway helps preserve fine-grained structural information. Combining M and R further enhances the discriminative ability of the feature space. When all four components (M, L, S, and R) are integrated, our complete RCDF achieves the highest performance, with Dice_{obj} = 74.89, AJI = 54.85, DQ = 70.65, SQ = 72.54, and PQ = 51.45, confirming the complementary roles of each submodule in achieving robust and accurate feature alignment.

7. Metrics Definition

Given that nuclei segmentation inherently targets instance-level predictions, the evaluation must assess how well individual nuclei are delineated rather than focusing solely on pixel-level accuracy. Therefore, two complementary object-centric metrics are employed: the *instance-level Dice coefficient* (D_{inst}) and the *Aggregated Jaccard Index* (AJI) [4]. The first metric, Dice_{obj}, quantifies the agreement between matched nuclei pairs and is defined as:

$$\begin{aligned}
 \text{Dice}_{\text{obj}}(Y, P) = & \frac{1}{2} \sum_{i=1}^{N_Y} \alpha_i D(Y_i, \tilde{P}(Y_i)) \\
 & + \frac{1}{2} \sum_{j=1}^{N_P} \beta_j D(\tilde{Y}(P_j), P_j),
 \end{aligned} \tag{1}$$

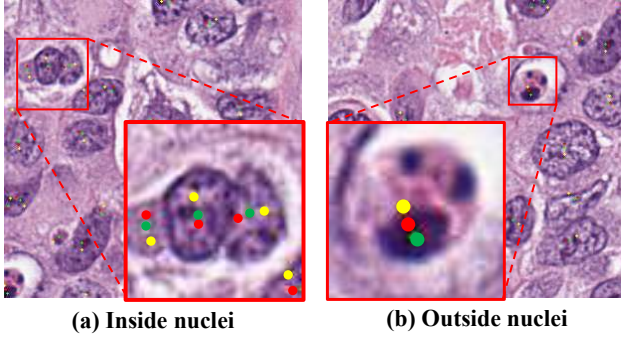


Figure 5. The figure illustrates simulated clinical annotation scenarios with different pixel offset distances. (a) shows all offset points located inside the nuclei, while (b) shows cases where some offset points fall outside the nuclei.

where $Y = \{Y_1, \dots, Y_{N_Y}\}$ and $P = \{P_1, \dots, P_{N_P}\}$ denote the sets of ground-truth and predicted nuclei, respectively. For each ground-truth nucleus Y_i , $\tilde{P}(Y_i)$ represents the predicted region that maximizes the overlap with Y_i , while $\tilde{Y}(P_j)$ indicates the ground-truth region having the greatest overlap with P_j . The weights α_i and β_j are proportional to the object areas. A one-to-one correspondence is established only if the overlapping region exceeds 50%. This formulation jointly evaluates prediction completeness and precision by averaging the Dice similarity across both annotation and prediction domains.

The *Aggregated Jaccard Index (AJI)* provides a complementary, global view of segmentation quality and is formulated as:

$$AJI = \frac{\sum_{i=1}^{N_Y} |Y_i \cap P^*(Y_i)|}{\sum_{i=1}^{N_Y} |Y_i \cup P^*(Y_i)| + \sum_{r \in \mathcal{R}} |P_r|}, \quad (2)$$

where $P^*(Y_i)$ is the predicted nucleus exhibiting the highest Jaccard similarity with Y_i , and \mathcal{R} is the set of unmatched predictions that do not correspond to any labeled instance. The numerator aggregates all intersection areas of valid matches, whereas the denominator accounts for both the union regions of matched pairs and the areas of remaining unmatched predictions. Consequently, J_{agg} serves as a holistic indicator that captures the overall consistency between predicted and ground-truth nuclei across the entire image.

8. Point Shift Visualization

Figure 5 illustrates the visualization of the simulated point-shift experiment designed to mimic real-world annotation deviations. Red dots represent randomly sampled points located 4 pixels away from the true nuclear centers (green dots), while yellow dots denote those 8 pixels away. In Figure 5 (a), all simulated points remain within the nu-

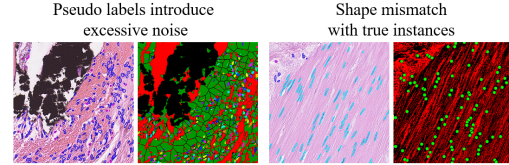


Figure 6. Visualization of the limitations of clustering-based pseudo labels on the CoNSeP dataset.

clear boundaries, which corresponds to the ideal annotation scenario. In contrast, Figure 5 (b) includes yellow points falling outside nuclei, simulating mislabeled cases commonly encountered in practice and testing the robustness of our method. As shown in the main paper (Figure 6), even when roughly one-quarter of the annotations are perturbed, our method still achieves strong performance. This robustness mainly stems from the complementary dual-representation design, which mitigates noise introduced by inaccurate annotations and reinforces nuclear features under both structural and semantic representations, leading to more stable instance delineation.

9. Limitation on the CoNSeP Dataset

As illustrated in Fig. 6, the CoNSeP dataset has relatively low image resolution, which makes the generation of clustering-based pseudo labels more prone to over-clustering and noise. In addition, pseudo labels produced by clustering methods tend to form compact and near-circular regions. However, many nuclei in the CoNSeP dataset exhibit elongated morphologies. This discrepancy between the pseudo-label shapes and the true nuclear structures can degrade the quality of the generated supervision and introduce additional noise during training. Such limitations are commonly encountered in weakly supervised nuclei instance segmentation when pseudo labels are derived from clustering-based strategies.

References

- [1] Simon Graham, Quoc Dang Vu, Shan E Ahmed Raza, Ayesha Azam, Yee Wah Tsang, Jin Tae Kwak, and Nasir Rajpoot. Hover-net: Simultaneous segmentation and classification of nuclei in multi-tissue histology images. *Medical image analysis*, 58:101563, 2019. 1
- [2] Neeraj Kumar, Ruchika Verma, Sanuj Sharma, Surabhi Bhargava, Abhishek Vahadane, and Amit Sethi. A dataset and a technique for generalized nuclear segmentation for computational pathology. *IEEE transactions on medical imaging*, 36(7):1550–1560, 2017. 1
- [3] Yi Lin, Zhiyong Qu, Hao Chen, Zhongke Gao, Yuexiang Li, Lili Xia, Kai Ma, Yefeng Zheng, and Kwang-Ting Cheng. Nuclei segmentation with point annotations from pathology images via self-supervised learning and co-training. *Medical Image Analysis*, 89:102933, 2023. 1
- [4] Hui Qu, Pengxiang Wu, Qiaoying Huang, Jingru Yi, Zhen-

- 287 nan Yan, Kang Li, Gregory M Riedlinger, Subhajyoti De,
288 Shaoting Zhang, and Dimitris N Metaxas. Weakly supervised
289 deep nuclei segmentation using partial points annotation in
290 histopathology images. *IEEE transactions on medical imag-*
291 *ing*, 39(11):3655–3666, 2020. 4
- 292 [5] Hui Qu, Jingru Yi, Qiaoying Huang, Pengxiang Wu, and Dim-
293 itris Metaxas. Nuclei segmentation using mixed points and
294 masks selected from uncertainty. In *2020 IEEE 17th Interna-*
295 *tional Symposium on Biomedical Imaging (ISBI)*, pages 973–
296 976. IEEE, 2020. 1
- 297 [6] Quoc Dang Vu, Simon Graham, Tahsin Kurc, Minh
298 Nguyen Nhat To, Muhammad Shaban, Talha Qaiser,
299 Navid Alemi Koohbanani, Syed Ali Khurram, Jayashree
300 Kalpathy-Cramer, Tianhao Zhao, et al. Methods for segmen-
301 tation and classification of digital microscopy tissue images.
302 *Frontiers in bioengineering and biotechnology*, 7:53, 2019. 1
- 303 [7] Kunzi Xie, Haonan Zhong, Jiamin Chang, Maurice Pagnucco,
304 and Yang Song. D2e2-net: Double deep edge enhancement for
305 weakly-supervised cell nuclei segmentation with incomplete
306 point annotations. In *2022 International Conference on Digi-*
307 *tal Image Computing: Techniques and Applications (DICTA)*,
308 pages 1–8. IEEE, 2022. 1