

## Highlights

### **Challenges in Deep Learning-Based Small Organ Segmentation: A Benchmarking Perspective for Medical Research with Limited Datasets**

Phongsakon Mark Konrad, Andrei-Alexandru Popa, Yaser Sabzehmeidani,  
Liang Zhong, Madhulika Tripathy, Andrei Constantinescu, Elisa A. Liehn,  
Serkan Ayvaz

- Benchmarking deep learning models on small medical datasets is highly unstable.
- Model rankings are dataset-specific and do not generalize across datasets.
- Statistical analysis reveals no significant difference between top models.
- Model rankings are sensitive to chosen cross-validation protocol and data splits.
- We advocate for uncertainty-aware evaluation in low-data clinical research scenarios.

# Challenges in Deep Learning-Based Small Organ Segmentation: A Benchmarking Perspective for Medical Research with Limited Datasets

Phongsakon Mark Konrad<sup>a</sup>, Andrei-Alexandru Popa<sup>b</sup>, Yaser Sabzehmeidani<sup>b</sup>,  
Liang Zhong<sup>c</sup>, Madhulika Tripathy<sup>c</sup>, Andrei Constantinescu<sup>e</sup>, Elisa A.  
Liehn<sup>d</sup>, Serkan Ayvaz<sup>a,\*</sup>

<sup>a</sup>*Centre for Industrial Software, University of Southern Denmark, Alsion  
2, Sønderborg, 6400, Denmark*

<sup>b</sup>*Centre for Industrial Mechanics, University of Southern Denmark, Alsion  
2, Sønderborg, 6400, Denmark*

<sup>c</sup>*Duke-NUS, 8 College Road, Singapore, 169857, Singapore*

<sup>d</sup>*National Heart Center Singapore, 5 Hospital Dr, Singapore, 169609, Singapore*

<sup>e</sup>*University of Medicine and Pharmacy Carol Davila Bucharest, Bulevardul Eroii Sanitari  
8, Bucharest, 050474, Romania*

---

## Abstract

Accurate segmentation of carotid artery structures in histopathological images is vital for cardiovascular disease research. This study systematically evaluates ten deep learning segmentation models including classical architectures, modern CNNs, a Vision Transformer, and foundation models, on a limited dataset of nine cardiovascular histology images. We conducted ablation studies on data augmentation, input resolution, and random seed stability to quantify sources of variance. Evaluation on an independent generalization dataset ( $N = 153$ ) under distribution shift reveals that foundation models maintain performance while classical architectures fail, and that rankings change sub-

---

\*Corresponding author

*Email addresses:* `phon@mimi.sdu.dk` (Phongsakon Mark Konrad), `andrei@sdu.dk` (Andrei-Alexandru Popa), `yasers@sdu.dk` (Yaser Sabzehmeidani), `zhong.liang@duke-nus.edu.sg` (Liang Zhong), `madhulika.tripathi@duke-nus.edu.sg` (Madhulika Tripathy), `andrei.constantinescu@umfcd.ro` (Andrei Constantinescu), `liehn.elisaanamarca@nhcs.com.sg` (Elisa A. Liehn), `seay@mimi.sdu.dk` (Serkan Ayvaz)

stantially between in-distribution and out-of-distribution settings. Training on the second dataset at varying sample sizes reveals dataset-specific ranking hierarchies confirming that model rankings are not generalizable across datasets. Despite rigorous Bayesian hyperparameter optimization, model performance remains highly sensitive to data splits. The bootstrap analysis reveals substantially overlapping confidence intervals among top models, with differences driven more by statistical noise than algorithmic superiority. This instability exposes limitations of standard benchmarking in low-data clinical settings and challenges assumptions that performance rankings reflect clinical utility. We advocate for uncertainty-aware evaluation in low-data clinical research scenarios from two perspectives. First, the scenario is not niche and is rather widely spread; and second, it enables pursuing or discontinuing research tracks with limited datasets from incipient stages of observations

*Keywords:* Medical Image Segmentation, Benchmark Instability, Limited Data, Statistical Stability, Histopathology, Foundation Models, Explainable AI

---

## 1 Introduction

Histological analysis is fundamental to cardiovascular pathology, providing detailed information for establishing diagnoses and tracking the morphological changes that define diseases like myocardial infarction and in-stent restenosis [1, 2, 3].

However, ensuring reproducibility of results across different laboratories remains a critical challenge. Although guidelines exist for standardizing instrument calibration and operation, documentation of microscopy protocols and data analysis procedures remains insufficient, compromising experimental reproducibility.

In the cardiovascular field, debates continue over the optimal methodologies for quantifying functional and morphological changes following surgical procedure in murine models, like atherosclerosis or myocardial infarction. Key questions revolve around how plaque size or infarct size should be measured, how it should be expressed (as a percentage based on a single section, averaged across multiple sections, and should focus solely on affected regions or include the entire organ). The absence of consensus has led to a heterogeneous body of published data, which complicates the translation of findings and impedes the development of standardized measurement protocols. As a result, researchers

20 must often reconstruct methodologies independently, slowing progress and  
21 reducing reproducibility in cardiovascular research

22 Computer-based systems and artificial intelligence (AI) can improve stan-  
23 dardization and consistency in cardiovascular research. While these tools  
24 are typically applied to large datasets, smaller datasets from animal studies  
25 receive less attention. The cardiovascular community increasingly seeks to  
26 minimize animal use [4] and obtain maximum information from minimal  
27 samples, even though current AI applications like 3D reconstructions do not  
28 fully support this approach. The variety of methods currently used raises  
29 concerns about the reliability of tissue analysis results. Therefore, efforts  
30 are being made to compare different methods on the same limited datasets  
31 to improve consistency and standardization across laboratories, ultimately  
32 strengthening the reproducibility of cardiovascular research.

33 As an example, the detailed characterization of vascular lesions from  
34 histological sections is especially critical in understanding conditions like  
35 carotid artery stenosis, a primary contributor to ischemic stroke [5]. However,  
36 the translation of this vital technique from research to widespread clinical  
37 application is present with challenges. The current gold standard, which relies  
38 on manual analysis of tissue by trained personnel, that is labor-intensive [6].  
39 This manual process not only risks damaging the fragile tissue samples but  
40 it introduces significant inter-observer variability, a fundamental barrier to  
41 producing reproducible science [7].

42 In recent years, deep learning models, from established Convolutional  
43 Neural Networks (CNNs) to massive, pre-trained Foundation Models (FMs),  
44 have shown promising results [8, 9, 10, 11]. However, these models are  
45 typically developed on large datasets, whereas specialized medical research  
46 often operates under data-scarce conditions [12]. This raises a fundamental  
47 methodological concern: when only a handful of patient samples are available,  
48 standard methods for selecting the “best” model may themselves be unreliable.

49 This study systematically investigates the performance of state-of-the-art  
50 deep learning models for carotid artery segmentation in the context of limited  
51 cardiovascular histopathological data. Our objective was to identify effective  
52 model architectures for this data-constrained task. We conducted a bench-  
53 marking study involving ten architectures: classical CNNs (FCN, SegNet),  
54 modern encoder-decoder networks (U-Net variants [13], DeepLabV3+ [14]),  
55 a lightweight transformer (SegFormer [15]), and foundation models (SAM  
56 [16], MedSAM [17], MedSAM+UNet [18]). Beyond performance compari-  
57 son, we conducted systematic ablation studies examining data augmentation

58 strategies (100 experiments across 10 presets), input resolution sensitivity  
59 (128–1024 pixels), and random seed stability (5 seeds per model) to quantify  
60 multiple sources of experimental variance.

61 The evaluation methodology followed widely adopted practices in ma-  
62 chine learning research, incorporating hyperparameter tuning via a thorough  
63 Bayesian optimization strategy and multiple cross-validation schemes [19].  
64 While such approaches aim to reduce variability and enhance reliability, they  
65 are not immune to the well-documented "illusion of control", a bias that  
66 can lead researchers to overestimate the stability of model performance in  
67 low-data regimes [20].

68 Beyond model selection, this study tests the hypothesis that standard  
69 benchmarking leaderboards are unreliable in low-data regimes, and that  
70 observed model rankings may reflect statistical noise and evaluation choices  
71 rather than true algorithmic differences. Using a dataset of  $N = 9$  images as a  
72 test case, we examine the limitations of common benchmarking practices when  
73 data is scarce, and propose an evaluation framework centered on statistical  
74 equivalence and practical efficiency

## 75 **2. Related Work**

### 76 *2.1. CNN-based Architectures for Medical Segmentation*

77 Convolutional Neural Networks (CNNs) have dominated medical image  
78 segmentation for over a decade. The U-Net architecture [13] remains the most  
79 widely adopted model, with numerous variants such as UNet++ [21] and  
80 R2U-Net [22] extending its encoder-decoder design. As noted by [23], these  
81 models have facilitated the translation of AI research to clinical application.  
82 DeepLabV3+ [14], with its Atrous Spatial Pyramid Pooling (ASPP) module  
83 for multi-scale feature extraction, provides an alternative approach. However,  
84 the performance of both U-Net and DeepLabV3+ remains fundamentally  
85 linked to the availability of large annotated datasets [24], a condition often  
86 unmet in specialized clinical research.

### 87 *2.2. Transformer-based Architectures in Vision and Medicine*

88 Transformer-based architectures leverage self-attention to model long-  
89 range dependencies, addressing the local receptive field limitations of CNNs  
90 [25]. SegFormer [15] combines a hierarchical Transformer encoder with a  
91 lightweight MLP decoder, offering an efficient alternative to CNN-based

92 segmentation. While such architectures provide broader representational  
93 capacity [26], they also introduce new validation challenges [27].

### 94 *2.3. Foundation Models and Their Adaptation for Medical Imaging*

95 The Segment Anything Model (SAM) [16], trained on over a billion masks,  
96 represents an important advance in general-purpose segmentation. Adapting  
97 SAM for medical imaging has been explored through domain-specific fine-  
98 tuning, as demonstrated by MedSAM [17], and through hybrid architectures  
99 that combine foundation model encoders with CNN-based decoders [18].  
100 However, the out-of-the-box performance of SAM on fine-grained medical  
101 tasks remains limited, necessitating careful adaptation strategies.

### 102 *2.4. Hyperparameter Optimization in Medical AI*

103 Model performance varies substantially with hyperparameter choices in-  
104 cluding learning rate, optimizer, scheduler, and loss function [28]. The result-  
105 ing high-dimensional search space makes exhaustive grid search infeasible,  
106 and random search provides no optimality guarantees [29]. Bayesian optimiza-  
107 tion addresses this by building probabilistic surrogate models to guide the  
108 search efficiently [30]. The challenge is compounded when Parameter-Efficient  
109 Fine-Tuning (PEFT) techniques such as LoRA [31] introduce additional  
110 hyperparameters that must themselves be tuned.

### 111 *2.5. The Benchmarking Reliability on Small Medical Datasets*

112 A broader methodological concern arises when comparing models in the  
113 low-data regimes common in medical research [32, 33, 34]. The practice of  
114 declaring a state-of-the-art model based on benchmark rankings is increasingly  
115 scrutinized [35], as a model’s rank can depend heavily on the evaluation  
116 protocol and data split [36]. LOOCV, while nearly unbiased, suffers from high  
117 variance [37], and extensive hyperparameter optimization risks overfitting to  
118 the validation set [38, 39], producing results that may not reproduce on new  
119 data [40]. Few studies address the stability of the evaluation process itself;  
120 our work addresses this gap by empirically quantifying how these concerns  
121 manifest in medical image segmentation.

## 122 **3. A Systematic Model Evaluation Framework**

123 To investigate the stability of deep learning benchmarks on small medical  
124 datasets, we designed a multi-stage evaluation framework. This process,

125 illustrated in Figure 1, was developed to ensure a fair and comprehensive  
 126 comparison across diverse model architectures by systematically controlling  
 127 for common sources of variability.

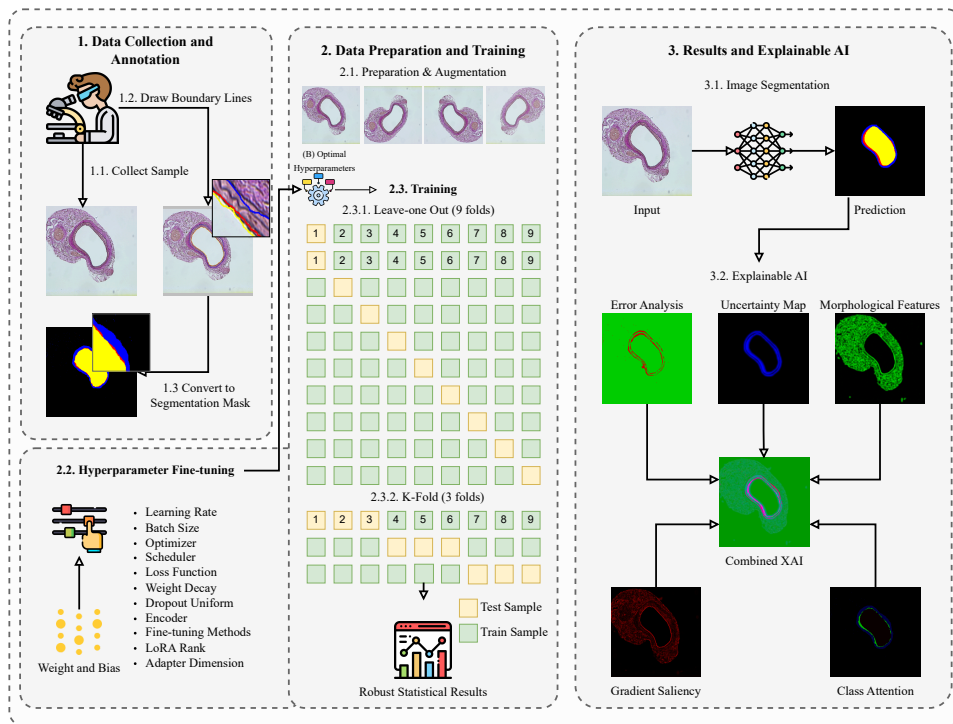


Figure 1: The systematic evaluation framework. Our process begins with (1) data collection and expert annotation. After the data preparation (2.1) Each model architecture then undergoes (2.2) extensive Bayesian hyperparameter optimization to find its optimal configuration and then use it for the training (2.3). (3.1) We then evaluate the segmentation performances, which are then dissected using our (3.2.) multi-modal XAI framework for qualitative insight.

128 The pipeline begins with data collection and expert annotation (1). Rather  
 129 than using arbitrary model configurations, each architecture undergoes system-  
 130 atic hyperparameter optimization (2.2.)—either Bayesian search or manual  
 131 tuning based on established practices—to ensure near-optimal performance.  
 132 This step is necessary for mitigating the confounding effect of suboptimal  
 133 configurations when comparing models [30, 28]. To generate statistically

134 sound results, we then subject each optimized model to multi-strategy cross-  
135 validation (2.3.1, 2.3.2) in combination with on-the-fly image augmentation  
136 (3.1), employing both Leave-One-Out (LOOCV) and 3-Fold splits. LOOCV,  
137 in particular, has been shown to provide robust statistical estimates for small  
138 datasets [41]. The quantitative results (4.1) are then dissected using our  
139 five-layer Explainable AI (XAI) framework (4.2) to identify potential sources  
140 of performance instability across data splits [42, 43].

### 141 3.1. Dataset and Clinical Context

142 The primary dataset (DS1) consists of nine high-resolution histological  
143 images of vascular tissue from the injured carotid artery of mice. Mice  
144 (C57/Bl6, LDL knockout male) undergo the carotid de-endothelialization,  
145 as described before by our group [44] (Permission no.231805 from IACUC  
146 at the Biology Resource Center at the Agency for Science, Technology and  
147 Research (A\*STAR), Singapore). The carotid artery was excised 2 weeks  
148 after de-endothelialization and embedded in paraffin. Serial sections (5  $\mu\text{m}$   
149 thick) were collected starting with the bifurcation between the internal carotid  
150 artery and external carotid artery. 10 sections, 5  $\mu\text{m}$  thick and 50  $\mu\text{m}$  apart  
151 were stained with orcein and photographed (Figure 3). These source images  
152 are stored in BMP format with variable native resolutions (e.g.,  $984 \times 792$   
153 pixels).

154 This small sample size ( $N = 9$ ) was intentionally chosen to reflect a clinical  
155 research scenario where data is scarce, requiring model robustness and data  
156 efficiency. The task is a four-class semantic segmentation problem designed  
157 to delineate key vascular structures: Lumen (class 1), Neointima (class 2),  
158 and Media (class 3), against the Background (class 0).

159 Ground truth annotations were independently performed by two domain  
160 experts (E.A. Liehn and A. Constantinescu) following the established an-  
161 notation protocol [44], which defines the delineation of lumen, neointima  
162 (intima), and media boundaries using the lamina interna and lamina externa  
163 as reference landmarks. The two annotators reached full agreement through  
164 the structured protocol; initial discrepancies were resolved by consensus before  
165 finalizing the ground truth. While a formal inter-annotator agreement score  
166 (e.g., inter-rater Dice) was not computed, the use of well-defined anatomical  
167 landmarks (lamina interna and lamina externa) as reference boundaries mini-  
168 mizes subjective variability, as demonstrated in prior work using this protocol  
169 [44].

170 Ground truth masks were generated from these expert annotations; initial  
171 color-coded JPG annotations were processed into four-class, integer-labeled  
172 PNG masks for model training and evaluation. For use in our pipeline, these  
173 high-resolution source images and masks are systematically preprocessed,  
174 including resizing to standardized input dimensions (e.g.,  $256 \times 256$  pixels for  
175 conventional models). An example is shown in Figure 2.

176 An independent generalization dataset (DS2) was collected from a sep-  
177 arate cohort of LDLR knockout mice (C57/Bl6 background) with carotid  
178 de-endothelialization performed as described above. Serial sections ( $5 \mu\text{m}$   
179 thick) from 22 carotid arteries were stained with orcein and photographed,  
180 yielding a total of  $N = 153$  images in TIF format. Three images were held out  
181 as a fixed independent test set; the remaining 150 were used for training. This  
182 dataset was used for two complementary experiments: (1) out-of-distribution  
183 generalization testing, where models trained on DS1 were evaluated directly  
184 on DS2 (Section 4.5); and (2) in-distribution training at varying sample sizes  
185 ( $N = 9, 25, 50, 100, 150$ ), where models were trained on nested DS2 subsets  
186 and evaluated on the three held-out DS2 images, to assess dataset-specific  
187 learning curves and ranking stability (Section 4.5.4). DS2 was not used during  
188 DS1-based hyperparameter selection. Ground truth annotations followed the  
189 same protocol as DS1, with segmentation of lumen, neointima, and media  
190 structures.

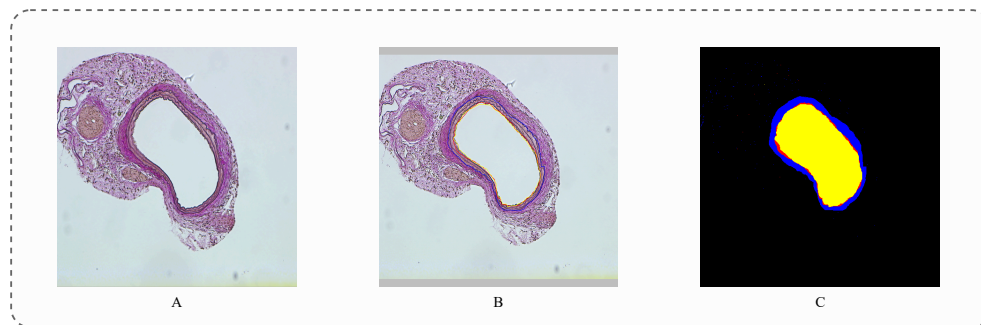


Figure 2: Example of data preparation. (A) Original histological image. (B) Expert line-art annotation. (C) Processed ground truth segmentation mask for **Lumen** (yellow), **Neointima** (red), and **Media** (blue).

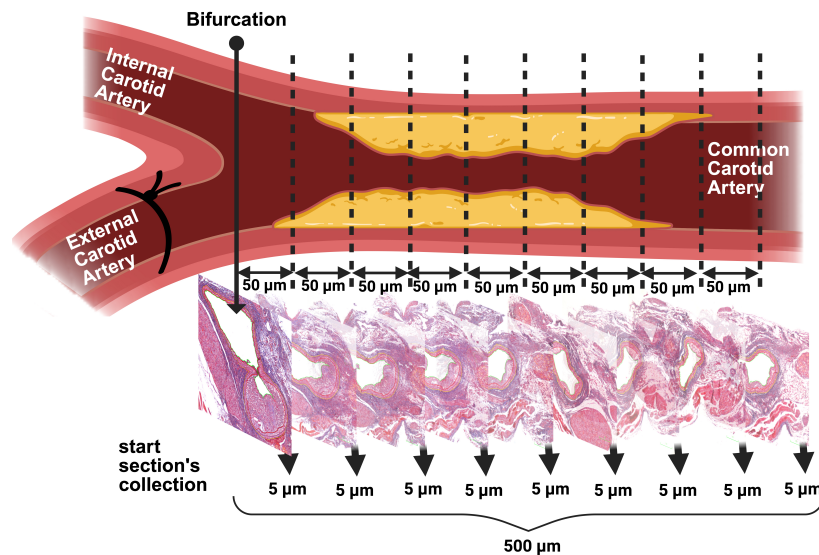


Figure 3: 5  $\mu\text{m}$  serial sections were collected starting from the bifurcation from the atherosclerotic injured common carotid artery, 50  $\mu\text{m}$  apart, as represented. Created in BioRender. Liehn, E. (2026) <https://BioRender.com/sm3xslo>.

### 191 3.2. Data Preprocessing and Augmentation

192 Our data preparation strategy employed a unified pipeline for all model  
 193 architectures to ensure a fair comparison. All images and their corresponding  
 194 masks were resized to a uniform dimension of  $256 \times 256$  pixels and normalized  
 195 to the range  $[0, 1]$ . For the foundation models (SAM and MedSAM), whose  
 196 Vision Transformer encoders (ViT-B, ViT-L) were originally designed for  
 197  $1024 \times 1024$  inputs, the upscaling to the encoder’s native resolution is handled  
 198 internally within each model’s forward pass via bilinear interpolation. The  
 199 output is subsequently downsampled back to the input resolution for loss  
 200 computation. This design ensures that all models receive identical input from  
 201 the data pipeline, ensuring consistent preprocessing across architectures.

202 Following this unified resizing, we employed an on-the-fly data augmenta-  
 203 tion strategy during training for all models to expand the limited dataset and  
 204 improve generalization [45, 46, 47]. This technique applied a series of random  
 205 transformations to each image-mask pair at every epoch, encompassing both  
 206 geometric and color-based changes. To account for variations in sample ori-  
 207 entation, geometric augmentations included random horizontal and vertical  
 208 flips, alongside affine transformations such as rotations ( $\pm 20^\circ$ ), scaling (from

209 80% to 120%), and shearing. These transforms were applied synchronously  
210 to both the image and its segmentation mask to maintain perfect spatial  
211 alignment. Color augmentations (jittering of brightness, contrast, saturation,  
212 and hue) were applied exclusively to input images to simulate variations in  
213 staining and lighting conditions.

214 To ensure a unbiased evaluation, the validation and test sets for all models  
215 were subjected only to their respective resizing and normalization steps,  
216 without any data augmentation.

### 217 *3.3. Model Architecture Selection*

218 We selected ten models spanning three architectural paradigms (classical  
219 and modern CNNs, a Vision Transformer, and Foundation Models) to provide  
220 a diverse testbed for evaluating benchmark stability.

#### 221 *3.3.1. Classical CNN Architectures*

222 FCN (134.4M parameters) uses a VGG-16 backbone with ImageNet-  
223 pretrained weights and serves as a baseline for dense prediction. SegNet  
224 (29.5M parameters) employs an encoder-decoder design that reuses max-  
225 pooling indices during upsampling. Both serve as historical baselines.

#### 226 *3.3.2. Modern CNN Architectures*

227 The standard UNet [13] (31.4M parameters), trained from scratch, serves  
228 as the canonical encoder-decoder baseline. To evaluate transfer learning, we  
229 compared eight pre-trained encoders from the ResNet [48], EfficientNet [49],  
230 DenseNet [50], and MobileNet [51] families; the best-performing configura-  
231 tion was UNet-ResNet34 (24.4M parameters). A Lightweight U-Net variant  
232 (7.8M parameters) with depthwise separable convolutions tests whether re-  
233 duced model capacity affects performance. DeepLabV3+ (26.7M parameters)  
234 captures multi-scale context through its ASPP module.

#### 235 *3.3.3. Modern Vision Transformer*

236 SegFormer employs a hierarchical Transformer encoder with a lightweight  
237 MLP decoder, providing a computationally efficient alternative for semantic  
238 segmentation (3.7M parameters).

#### 239 *3.3.4. Foundation Models and Hybrid Architectures*

240 We included SAM and its medical derivative MedSAM to assess the  
241 impact of large-scale pre-training. Since SAM requires manual prompts for

242 segmentation, which is impractical for automated benchmarking [52, 53], we  
243 implemented a learned auto-prompting module: a lightweight network uses  
244 multi-head cross-attention to fuse SAM’s image embeddings with learnable  
245 prompt tokens, enabling fully automated end-to-end training. MedSAM was  
246 fine-tuned using Norm Tuning, a parameter-efficient strategy that updates  
247 only the normalization layers of the frozen encoder alongside the mask decoder.  
248 We also included a MedSAM-UNet hybrid that feeds the frozen MedSAM  
249 encoder features into a U-Net-style decoder, combining the context modeling  
250 of a Transformer with the spatial precision of a CNN decoder.

### 251 3.4. Implementation Reproducibility

252 To ensure reproducibility and enable fair comparison across experimental  
253 conditions, we document the complete computational environment. All  
254 experiments were conducted using NVIDIA Tesla V100 GPUs (32GB VRAM)  
255 with PyTorch 2.1+ and CUDA 12.1. Automatic Mixed Precision (AMP)  
256 training was enabled to improve computational efficiency.

#### 257 3.4.1. Randomness Control

258 Our implementation sets seeds across all relevant random number gen-  
259 erators using a unified seed value ( $seed = 42$ ), including Python’s random  
260 module, NumPy’s random state, and PyTorch CPU/CUDA operations. To  
261 ensure deterministic behavior, we configured PyTorch’s cuDNN backend with  
262 deterministic mode enabled and benchmark mode disabled.

#### 263 3.4.2. Limitations of Single-Seed Evaluation

264 A methodological limitation is the use of a single random seed for all  
265 experiments. While this ensures reproducibility of our specific results, it  
266 does not capture the variance from different random initializations. This  
267 limitation reinforces our central finding: even with controlled randomness,  
268 the small dataset size introduces substantial variance through the data splits  
269 themselves.

### 270 3.5. Hyperparameter Optimization

271 Hyperparameters significantly impact model performance, particularly  
272 with limited data [28]. To account for this, we used Bayesian optimization  
273 (explained in Appendix A) managed by Optuna<sup>1</sup> [54], a Bayesian hyperpa-

---

<sup>1</sup><https://optuna.org/>

274 parameter optimization framework using the Tree-structured Parzen Estimator  
 275 (TPE) sampler. Each of the ten architectures was independently optimized  
 276 over 100 Bayesian trials, totaling 1,000 runs. Our extensive search space is  
 277 detailed in Table 1. As Figure 4 illustrates, this approach is crucial: without  
 278 it, a strong architecture might be unfairly dismissed due to a single suboptimal  
 279 configuration [55]. This process ensures each architecture is compared at its  
 280 near-optimal configuration, which is detailed in Table 2.

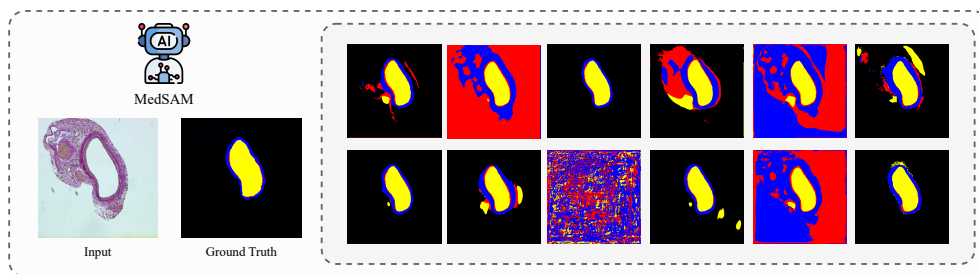


Figure 4: The impact of hyperparameter selection. All twelve segmentation results were generated by the same MedSAM architecture. The wide variation in output quality is due solely to different hyperparameter configurations. This illustrates that without a rigorous search, one could incorrectly dismiss a capable model based on a single, suboptimal trial.

### 281 3.6. Transfer Learning and Weight Initialization Strategy

282 Fair model comparison requires consistent initialization strategies. Table  
 283 3 summarizes the pretraining source, weight initialization, and fine-tuning  
 284 approach for each architecture.

#### 285 3.6.1. CNN and Transformer Models: Full Fine-tuning

286 For the encoder-based CNN architectures (UNet-ResNet34, DeepLabV3+)  
 287 and SegFormer, we employed full fine-tuning with ImageNet-pretrained en-  
 288 coders via the `segmentation-models-pytorch` library. The encoder weights  
 289 were initialized with ImageNet weights but were not frozen during training,  
 290 allowing the encoder to adapt its learned representations to the histopatho-  
 291 logical domain. The decoders were randomly initialized. The remaining  
 292 architectures (UNet, Lightweight U-Net, FCN, SegNet) were trained with  
 293 random initialization of all weights.

Table 1: Hyperparameter Search Space Summary.

Parameter	Search Strategy	Range / Values
<i>General Parameters</i>		
Learning Rate	Log-uniform	[1e-5, 1e-2]
Batch Size	Categorical	[2, 4, 8, 16, 32]
Optimizer	Categorical	[Adam, AdamW, SGD, RMSprop, Nadam]
Scheduler	Categorical	[Cosine, Polynomial, Step, Warmup Cosine, OneCycle]
Loss Function	Categorical	[Focal-Dice, Focal-Tversky, Unified Focal]
Weight Decay	Log-uniform	[1e-6, 1e-2]
Dropout	Uniform	[0.0, 0.3]
<i>Architecture-Specific Parameters</i>		
Encoder (UNet)	Categorical	[ResNet18/34/50, EfficientNet-B0/B1/B2, DenseNet121, MobileNetV2]
SAM Model Type	Categorical	[ViT-B, ViT-L]
Fine-tuning Method	Categorical	[SVD, LoRA, Norm Tuning]
LoRA Rank	Categorical	[8, 16, 32]
Adapter Dimension	Categorical	[128, 256, 512]

294 *3.6.2. Foundation Models: Parameter-Efficient Fine-tuning*

295 For the foundation models (SAM, MedSAM, MedSAM+UNet), we adopted  
 296 a parameter-efficient fine-tuning (PEFT) strategy. The image encoders were  
 297 initialized from official SAM checkpoints and were frozen by default to reduce  
 298 computational requirements and prevent overfitting. All three foundation  
 299 models use Norm Tuning as their PEFT method, which enables gradient  
 300 computation only for the normalization layers within the frozen encoder while  
 301 training the mask decoder. For SAM, this was identified as the optimal  
 302 method via Bayesian search, combined with a learned auto-prompting module.  
 303 For MedSAM, Norm Tuning allows the normalization layers and mask decoder  
 304 to adapt to the histopathological domain. For MedSAM+UNet, the SAM  
 305 encoder normalization layers are tuned alongside the custom UNet-style

Table 2: Optimal Hyperparameters for Each Model. All ten models were independently optimized via Bayesian hyperparameter search (100 trials each). All models use Focal-Dice loss and train for 200 epochs.

	FCN	SegNet	LW-U <sup>d</sup>	UNet	U-R34 <sup>c</sup>	DLV3+ <sup>a</sup>	SF <sup>b</sup>	SAM	MS <sup>h</sup>	M-U <sup>e</sup>
<i>Architecture</i>										
Backbone	VGG-16	VGG	Depth-sep	Standard	RN34 <sup>f</sup>	RN50	MiT-B0	ViT-B	ViT-L	ViT-B+U
Pretrained	IN <sup>f</sup>	—	—	—	IN	IN	IN	SAM	SAM	SAM
<i>Optimization</i>										
LR	1e-4	1e-4	1e-4	1e-4	1e-4	5e-5	6e-5	6e-4	8.2e-4	8.7e-4
Optimizer	AdamW	AdamW	AdamW	AdamW	Adam	Adam	Adam	AdamW	Adam	Adam
Scheduler	1Cyc <sup>§</sup>	1Cyc	1Cyc	1Cyc	1Cyc	1Cyc	Cos	WCos <sup>§</sup>	WCos	WCos
Batch	8	8	32	8	8	28	8	4	2	4
WD	1e-4	1e-4	1e-4	1e-4	1e-4	1e-4	1e-4	1e-3	6.7e-4	8.1e-5
<i>Regularization &amp; PEFT</i>										
Dropout	0.5	0.5	0.1	0.1	0.1	—	—	—	—	0.0
Grad. Clip	—	—	—	—	—	—	—	2.0	1.0	1.0
PEFT	—	—	—	—	—	—	—	NT <sup>§</sup>	NT	NT
Adapter Dim.	—	—	—	—	—	—	—	512	512	—

<sup>a</sup> DLV3+: DeepLabV3+; <sup>b</sup> SF: SegFormer; <sup>c</sup> U-R34: UNet-ResNet34; <sup>d</sup> LW-U: Lightweight UNet; <sup>e</sup> M-U: MedSAM+UNet. <sup>h</sup> MS: MedSAM. <sup>f</sup> RN: ResNet; IN: ImageNet. <sup>§</sup> 1Cyc: OneCycle; WCos: Warmup Cosine; Cos: Cosine; NT: Norm Tuning; WD: Weight Decay; LR: Learning Rate.

Table 3: Transfer Learning and Weight Initialization Strategy. All decoders are randomly initialized. Params indicates total model parameters.

Model	Pretraining	Encoder Frozen?	PEFT	Params (M)	Trainable Scope
<i>Classical Architectures</i>					
FCN	ImageNet (VGG-16)	No	—	134.4	Full network
SegNet	None	N/A	—	29.5	Full network
<i>CNN-based Architectures</i>					
UNet	None	N/A	—	31.4	Full network
LW-UNet	None	N/A	—	7.8	Full network
UNet-ResNet34	ImageNet	No	—	24.4	Encoder + decoder
DeepLabV3+	ImageNet	No	—	26.7	Encoder + decoder + ASPP
<i>Transformer-based Architecture</i>					
SegFormer	ImageNet	No	—	3.7	MiT encoder + MLP decoder
<i>Foundation Models</i>					
SAM	SA-1B	Yes	Norm Tuning	97.3	Norm layers + decoder + auto-prompt
MedSAM	SA-1B + Medical	Yes	Norm Tuning	95.8	Norm layers + mask decoder
MedSAM+UNet	SA-1B + Medical	Yes	Norm Tuning	95.7	Norm layers + UNet decoder

306 decoder.

307 *3.7. Evaluation Metrics and Cross-Validation*

308 Our primary evaluation metric is the macro-averaged Dice Similarity  
 309 Coefficient (DSC) over the three tissue classes, chosen for its robustness to  
 310 class imbalance [56]. We also report Intersection over Union (IoU) [57]. To  
 311 assess ranking stability, we employ two cross-validation strategies: Leave-  
 312 One-Out (LOOCV), a low-bias, high-variance estimator, and 3-Fold CV, a  
 313 higher-bias, lower-variance estimator [58].

314 For a given ground truth segmentation mask  $\mathbf{Y} \in \{0, 1, 2, 3\}^{H \times W}$  and  
 315 its corresponding predicted segmentation  $\hat{\mathbf{Y}} \in \{0, 1, 2, 3\}^{H \times W}$ , with classes  
 316 defined as  $\{0 : \text{Background}, 1 : \text{Lumen}, 2 : \text{Neointima}, 3 : \text{Media}\}$ , we define  
 317 the fundamental confusion matrix elements for each class  $c$ . The indicator  
 318 function  $\mathbf{1}[\cdot]$  evaluates to 1 if the condition is true and 0 otherwise.

$$TP_c = \sum_{i,j} \mathbf{1}[\mathbf{Y}(i, j) = c \wedge \hat{\mathbf{Y}}(i, j) = c] \quad (1)$$

$$FP_c = \sum_{i,j} \mathbf{1}[\mathbf{Y}(i, j) \neq c \wedge \hat{\mathbf{Y}}(i, j) = c] \quad (2)$$

$$FN_c = \sum_{i,j} \mathbf{1}[\mathbf{Y}(i, j) = c \wedge \hat{\mathbf{Y}}(i, j) \neq c] \quad (3)$$

319 where  $TP_c$  (True Positives) counts pixels correctly classified as class  $c$ ,  $FP_c$   
 320 (False Positives) counts pixels incorrectly predicted as class  $c$ , and  $FN_c$  (False  
 321 Negatives) counts pixels of class  $c$  misclassified as other classes.

The Dice Similarity Coefficient ( $DSC_c$ ) quantifies the spatial overlap  
 between the predicted and ground truth segmentations for a given class  $c$ :

$$DSC_c = \frac{2 \cdot TP_c}{2 \cdot TP_c + FP_c + FN_c}$$

To ensure our optimization targets clinically relevant tissues, we compute  
 a macro-averaged Dice coefficient ( $DSC_{\text{macro}}$ ) by excluding the background  
 class ( $c = 0$ ):

$$DSC_{\text{macro}} = \frac{1}{3} \sum_{c=1}^3 DSC_c$$

Intersection over Union ( $IoU_c$ ), also known as the Jaccard index, provides  
 a complementary measure of segmentation quality for class  $c$ :

$$IoU_c = \frac{TP_c}{TP_c + FP_c + FN_c}$$

The relationship between DSC and IoU is direct:

$$DSC_c = \frac{2 \cdot IoU_c}{1 + IoU_c}$$

322 Our loss function, identified as optimal through our hyperparameter search,  
 323 is a composite of Focal and Dice loss designed to balance pixel-level accuracy  
 324 and region-level overlap [59].

### 325 3.7.1. Focal Loss

To mitigate the impact of severe class imbalance, we employ Focal Loss ( $\mathcal{L}_{\text{focal}}$ ). This loss function down-weights easy examples and focuses training on hard, misclassified examples:

$$\mathcal{L}_{\text{focal}} = -\alpha(1 - p_t)^\gamma \log(p_t)$$

326 Here,  $p_t$  represents the model’s estimated probability for the ground truth  
 327 class. We set  $\alpha = 0.8$  for class balancing and  $\gamma = 2.0$  to control the focusing  
 328 strength on hard examples.

### 329 3.7.2. Dice Loss

330 The Dice Loss ( $\mathcal{L}_{\text{dice}}$ ) directly optimizes our primary evaluation metric,  
 331 promoting better overlap between predicted and ground truth masks. Its  
 332 differentiable form is given by:

$$\mathcal{L}_{\text{dice}} = 1 - \frac{1}{C} \sum_{c=1}^C \frac{2 \sum_{i,j} \mathbf{Y}_c(i,j) \cdot \hat{\mathbf{Y}}_c(i,j) + \epsilon}{\sum_{i,j} \mathbf{Y}_c(i,j) + \sum_{i,j} \hat{\mathbf{Y}}_c(i,j) + \epsilon} \quad (4)$$

333 In this formulation,  $\hat{\mathbf{Y}}_c$  denotes the predicted probabilities for class  $c$ ,  $\mathbf{Y}_c$   
 334 is the one-hot encoded ground truth for class  $c$ , and  $\epsilon = 10^{-6}$  is a small constant  
 335 added to prevent division by zero.  $C$  is the total number of classes.

### 336 3.7.3. Composite Loss Function

Our primary and most effective loss function, the Composite Loss Function ( $\mathcal{L}_{\text{focal\_dice}}$ ), strategically combines Focal Loss and Dice Loss with adaptive weighting:

$$\mathcal{L}_{\text{focal\_dice}} = \lambda_f \mathcal{L}_{\text{focal}} + \lambda_d \mathcal{L}_{\text{dice}}$$

337 Here,  $\lambda_f = 0.5$  and  $\lambda_d = 0.5$  provide a balanced optimization of both pixel-  
 338 level accuracy and region-level overlap, leading to superior segmentation  
 339 results.

340 *3.8. Multi-Modal Explainable AI (XAI) Framework*

341 Interpretability supports clinical adoption of automated segmentation.  
342 We present a framework that integrates five XAI modalities (Figure 5) to  
343 provide a multi-faceted diagnostic view. We apply this framework not only  
344 to explain individual predictions but also to diagnose sources of statistical  
345 instability across different data splits. An explainability framework is only  
346 useful if its outputs map directly to clinical needs. We designed our five layers  
347 to address distinct, practical questions that arise during pathological review,  
348 from assessing diagnostic accuracy to building confidence in the model’s  
349 predictions. Table 4 provides a systematic validation of how each component  
350 serves critical clinical functions.

Mathematically, our framework generates a composite explanation map  $\mathbf{E}$  from an input image  $\mathbf{I}$ , its ground truth  $\mathbf{Y}$ , and the model’s prediction  $\hat{\mathbf{Y}}$ . This composite map is constructed as a weighted sum of five individual, distinct explanation layers  $\mathbf{E}_i$ :

$$\mathbf{E} = \sum_{i=1}^5 \alpha_i \mathbf{E}_i$$

351 The coefficients  $\alpha_i$  are adaptive and can be adjusted to emphasize different  
352 explanatory aspects based on the specific clinical question or the user’s focus.  
353 The individual layers are detailed below:

354 *3.8.1. Layer 1: Error Analysis*

This foundational layer provides direct visual feedback on the correctness of the model’s prediction. It generates a binary error map  $\mathbf{E}_1$  that immediately highlights regions of agreement and disagreement between the predicted mask  $\hat{\mathbf{Y}}$  and the ground truth  $\mathbf{Y}$ .

$$\mathbf{E}_1(i, j) = \begin{cases} \blacksquare \text{ (Green, RGB: 0, 204, 0)} & \text{if } \hat{\mathbf{Y}}(i, j) = \mathbf{Y}(i, j) \\ \blacksquare \text{ (Red, RGB: 204, 0, 0)} & \text{otherwise} \end{cases}$$

355 Here, **green** represents correctly classified pixels where the prediction matches  
356 the ground truth, while **red** highlights erroneous pixels where the model’s  
357 prediction diverges from the expert annotation.

358 *3.8.2. Layer 2: Uncertainty Estimation*

This layer quantifies the model’s confidence in its pixel-wise predictions. Uncertainty is often highest at critical tissue boundaries, indicating areas

where the model is less certain about its classification. We compute an uncertainty map  $\mathbf{E}_2$  using the pixel-wise entropy  $\mathcal{H}$  of the model’s class probabilities  $\mathbf{P}$ .

$$\mathbf{E}_2(i, j) = \mathcal{H}(\mathbf{P}(i, j)) = - \sum_c P_c(i, j) \log_2 P_c(i, j)$$

359 *3.8.3. Layer 3: Morphological Analysis*

360 To ground the explanation in clinically relevant anatomical and structural  
 361 features, this layer  $\mathbf{E}_3$  captures tissue-specific morphological characteristics by  
 362 integrating various image processing techniques, such as texture descriptors  
 363 and structural filters.

364 *3.8.4. Layer 4: Class-wise Attention*

This layer  $\mathbf{E}_4$  adaptively highlights regions corresponding to specific classes where the model exhibits suboptimal performance. The contribution of each class mask  $\mathbf{M}_c$  to the overall attention map is weighted by its performance deficit, calculated as  $(1 - DSC_c)$ . This approach focuses the explanation on the most problematic tissue types.

$$\mathbf{E}_4(i, j) = \sum_{c=1}^3 (1 - DSC_c) \cdot \mathbf{M}_c(i, j)$$

365 *3.8.5. Layer 5: Gradient-based Saliency*

The final layer  $\mathbf{E}_5$  emphasizes the fine-grained structural edges that are often critical for accurate histological interpretation. This is achieved by computing the magnitude of the image gradient, highlighting areas of rapid intensity change.

$$\mathbf{E}_5(i, j) = \sqrt{\left(\frac{\partial \mathbf{I}}{\partial i}\right)^2 + \left(\frac{\partial \mathbf{I}}{\partial j}\right)^2}$$

366 This multi-modal framework provides clinicians with a detailed view of model  
 367 behavior, supporting informed diagnostic decisions.

368 *3.9. Statistical Analysis*

369 Given the small sample size of our dataset ( $N = 9$ ), we cannot assume a  
 370 normal distribution for performance scores. Therefore, our statistical analysis  
 371 relies primarily on non-parametric methods to compare model performance

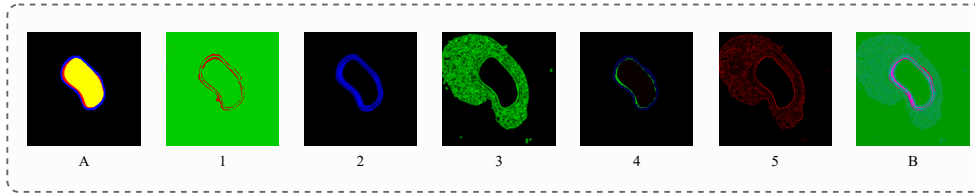


Figure 5: The five-layer XAI framework. (A) A model’s prediction is dissected into (1) Error Analysis, (2) Uncertainty, (3) Morphology, (4) Attention, and (5) Saliency. (B) These are synthesized into an integrated explanation, providing a multi-faceted view of the model’s behavior.

Table 4: Clinical Validation Framework for XAI Components. Column headers (1-5, B) refer to the components detailed in Figure 5: (1) Error, (2) Uncertainty, (3) Morphology, (4) Attention, (5) Saliency, and (B) Integrated.

Clinical Need	(1)	(2)	(3)	(4)	(5)	(B)
Diagnostic Accuracy	✓	✓	✓	✓	✓	✓
Boundary Delineation	×	✓	✓	×	✓	✓
Tissue Characterization	×	×	✓	×	×	✓
Quality Assurance	✓	✓	×	✓	×	✓
Educational Value	✓	✓	✓	✓	✓	✓
Clinical Confidence	×	✓	✓	✓	×	✓

372 and assess the stability of the results. All statistical analyses were con-  
 373 ducted using Python, leveraging the *pandas* library for data manipulation  
 374 and *Matplotlib* with *Seaborn* for visualization.

375 To quantify the uncertainty around mean performance scores without  
 376 assuming a normal distribution, we calculated 95% bootstrap confidence  
 377 intervals [60]. This was achieved by resampling the cross-validation fold scores  
 378 with replacement 10,000 times for each model to generate a distribution of  
 379 the mean, from which the percentile-based confidence intervals were derived.  
 380 The core resampling was implemented using the `resample` function from  
 381 *scikit-learn*<sup>2</sup> [61].

382 To formally compare the models against each other across all cross-

<sup>2</sup><https://scikit-learn.org/>

383 validation folds, we employed the Friedman test [62]. This non-parametric  
384 test is highly recommended for comparing multiple classifiers over multiple  
385 datasets (or, in our case, folds) [63]. We then used the Nemenyi test as a  
386 post-hoc analysis to identify which specific models have statistically different  
387 average ranks. The results are visualized using Critical Difference (CD) plots.  
388 This entire analysis was performed using functions from the *SciPy*<sup>3</sup> library  
389 [64].

390 Finally, to assess the practical magnitude of the performance differences  
391 between model pairs, we calculated Cohen’s d [65]. This provides a stan-  
392 dardized measure of the difference between two means, helping to distinguish  
393 between statistical significance and practical importance.

## 394 4. Results

395 This section details the empirical findings of our study. We first present  
396 quantitative performance under two cross-validation protocols, then examine  
397 the stability of these results through statistical analysis.

### 398 4.1. Quantitative Performance Across Validation Protocols

399 To establish a performance baseline, all optimally-configured models were  
400 evaluated using both Leave-One-Out (LOOCV) and 3-Fold cross-validation  
401 (CV).

402 Under the LOOCV protocol, which provides a nearly unbiased but high-  
403 variance estimate of performance, MedSAM emerged as the apparent leader.  
404 As shown in Table 5, MedSAM achieved a macro-averaged Dice Similarity  
405 Coefficient (DSC) of 0.694, followed by SegFormer at 0.576. With the ex-  
406 panded set of ten models, the classical architectures (FCN, SegNet) performed  
407 substantially worse, achieving DSC scores below 0.15.

408 Under the 3-Fold CV protocol, which has lower variance but potentially  
409 higher bias (Table 6), MedSAM retained the top position (0.643 DSC),  
410 followed by SegFormer (0.501 DSC). While MedSAM leads in both protocols,  
411 substantial rank instability exists among other models—as visualized in Figure  
412 6. While MedSAM leads in both protocols, the statistical robustness of this  
413 ranking is examined in the following sections.

---

<sup>3</sup><https://scipy.org/>

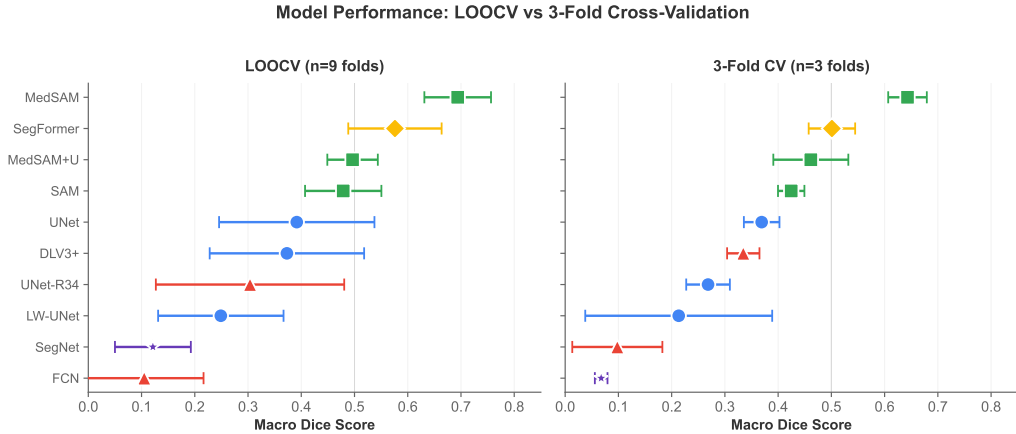


Figure 6: Performance Comparison Across Validation Protocols. Side-by-side comparison of model rankings under LOOCV (left) and 3-Fold CV (right). While MedSAM leads under both protocols, substantial rank instability exists among other models—the crossing lines highlight how relative positions change between protocols.

Table 5: LOOCV Performance Summary (10 Models). Results are shown as Mean  $\pm$  Std Dev. Models ranked by Macro Dice score; MedSAM achieves the highest mean.

Model	Macro Dice $\uparrow$	Macro IoU $\uparrow$
<i>Foundation Models &amp; Transformers</i>		
<b>MedSAM</b>	<b>0.694 <math>\pm</math> 0.063</b>	<b>0.592 <math>\pm</math> 0.058</b>
SegFormer	0.576 $\pm$ 0.088	0.474 $\pm$ 0.074
SAM	0.496 $\pm$ 0.047	0.412 $\pm$ 0.035
MedSAM+UNet	0.479 $\pm$ 0.072	0.402 $\pm$ 0.068
<i>Modern CNNs</i>		
UNet	0.392 $\pm$ 0.146	0.319 $\pm$ 0.122
UNet-ResNet34	0.373 $\pm$ 0.145	0.297 $\pm$ 0.132
DeepLabV3+	0.304 $\pm$ 0.177	0.240 $\pm$ 0.150
Lightweight UNet	0.249 $\pm$ 0.118	0.200 $\pm$ 0.102
<i>Classical Architectures</i>		
FCN	0.121 $\pm$ 0.071	0.075 $\pm$ 0.055
SegNet	0.105 $\pm$ 0.111	0.076 $\pm$ 0.088

#### 4.2. Analysis of Benchmark Instability

Further analysis indicates that no single model can be reliably identified as superior, as rankings depend on the specifics of the evaluation protocol. We identified three primary sources of this instability.

Table 6: 3-Fold Cross-Validation Performance Summary (10 Models). Results are shown as Mean  $\pm$  Std Dev across 3 folds. Note the ranking changes compared to LOOCV.

Model	Macro Dice $\uparrow$	Macro IoU $\uparrow$
<i>Foundation Models &amp; Transformers</i>		
<b>MedSAM</b>	<b>0.643 <math>\pm</math> 0.036</b>	<b>0.542 <math>\pm</math> 0.032</b>
SegFormer	0.501 $\pm$ 0.044	0.415 $\pm$ 0.040
MedSAM+UNet	0.461 $\pm$ 0.070	0.388 $\pm$ 0.072
SAM	0.425 $\pm$ 0.025	0.351 $\pm$ 0.029
<i>Modern CNNs</i>		
UNet	0.369 $\pm$ 0.033	0.289 $\pm$ 0.034
DeepLabV3+	0.335 $\pm$ 0.030	0.268 $\pm$ 0.033
UNet-ResNet34	0.268 $\pm$ 0.041	0.201 $\pm$ 0.041
Lightweight UNet	0.213 $\pm$ 0.176	0.169 $\pm$ 0.149
<i>Classical Architectures</i>		
SegNet	0.098 $\pm$ 0.085	0.069 $\pm$ 0.069
FCN	0.068 $\pm$ 0.012	0.036 $\pm$ 0.006

#### 4.2.1. Sensitivity of Model Rankings to Evaluation Protocol

The most direct evidence of instability is that the top-ranked model changes depending on the evaluation context. Table 7 summarizes the best-performing model for each metric across three contexts: the 80/20 train-test split using optimized hyperparameters, the mean LOOCV result, and the mean 3-Fold CV result. While MedSAM wins across all evaluation protocols in this table, it is notable that during hyperparameter optimization, UNet achieved the highest validation Dice score (0.630) compared to MedSAM (0.569)—yet this ranking reverses on held-out test data. This exemplifies how validation-based model selection can be misleading.

Furthermore, the substantial range between highest and lowest scores across all 10 models (e.g., 0.007–0.781 for LOOCV Dice, representing a 0.774 range) highlights the substantial performance variability within each protocol.

#### 4.2.2. Mean Scores Obscure High Fold-to-Fold Variance

Average performance metrics, while useful, can mask substantial variation across individual data splits. This is visualized in Figure 7, which plots the performance rank of each model on each fold of the cross-validation. The left panel, showing 3-Fold CV, exhibits relatively consistent rankings. By contrast, the right panel for LOOCV reveals substantial rank variability. For example, while MedSAM maintains a stable 1st-place rank across all folds, models such as UNet-ResNet34 and DeepLabV3+ exhibit rank swings of up

Table 7: Analysis of Performance Instability Across Evaluation Contexts (10 Models). The Range ( $\Delta$ ) column shows performance spread across all 10 models, quantifying instability magnitude.

<b>Metric</b>	<b>Winner</b>	<b>Highest</b>	<b>Lowest</b>	<b>Range (<math>\Delta</math>)</b>
<b>Best Single Run (80/20 Split with Optimized Hyperparameters)</b>				
<i>Macro-Averaged Metrics (Overall Performance)</i>				
Dice (Macro)	MedSAM	0.609	0.057	0.552
IoU (Macro)	MedSAM	0.513	0.030	0.483
<i>Class-Specific Metrics (Tissue-Level Performance)</i>				
Dice (Lumen)	MedSAM	0.941	0.125	0.816
Dice (Neointima)	MedSAM	0.224	0.000	0.224
Dice (Media)	MedSAM	0.662	0.040	0.622
IoU (Lumen)	MedSAM	0.903	0.067	0.836
IoU (Neointima)	MedSAM	0.129	0.000	0.129
IoU (Media)	MedSAM	0.507	0.020	0.487
<b>Leave-One-Out CV (LOOCV)</b>				
<i>Macro-Averaged Metrics (Overall Performance)</i>				
Dice (Macro)	MedSAM	0.781	0.007	0.774
IoU (Macro)	MedSAM	0.676	0.003	0.673
<i>Class-Specific Metrics (Tissue-Level Performance)</i>				
Dice (Lumen)	MedSAM	0.974	0.002	0.972
Dice (Neointima)	MedSAM	0.715	0.000	0.715
Dice (Media)	MedSAM	0.788	0.000	0.788
IoU (Lumen)	MedSAM	0.950	0.001	0.949
IoU (Neointima)	MedSAM	0.582	0.000	0.582
IoU (Media)	MedSAM	0.656	0.000	0.656
<b>3-Fold CV (K-Fold)</b>				
<i>Macro-Averaged Metrics (Overall Performance)</i>				
Dice (Macro)	MedSAM	0.688	0.022	0.667
IoU (Macro)	MedSAM	0.584	0.011	0.573
<i>Class-Specific Metrics (Tissue-Level Performance)</i>				
Dice (Lumen)	MedSAM	0.958	0.001	0.957
Dice (Neointima)	MedSAM	0.451	0.002	0.449
Dice (Media)	MedSAM	0.712	0.054	0.658
IoU (Lumen)	MedSAM	0.921	0.000	0.920
IoU (Neointima)	MedSAM	0.325	0.001	0.323
IoU (Media)	MedSAM	0.570	0.028	0.542

439 to 7 positions between folds. This indicates that a high average score does not  
 440 guarantee reliable performance on any single, unseen sample, an important  
 441 concern for clinical applications.

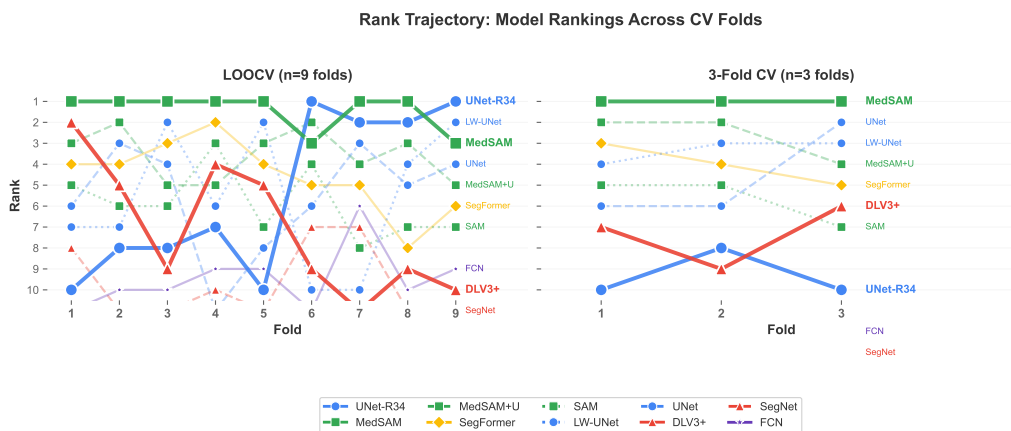


Figure 7: Model Rank Stability in 3-Fold vs. Leave-one-out cross-validation (LOOCV). The plots track the performance rank of each model (1=Best) across the validation folds for 3-Fold CV (left) and LOOCV (right). The rank trajectories reveal substantial fold-to-fold volatility, particularly under LOOCV.

#### 4.2.3. Statistical Analysis of Performance Differences

442 To rigorously test the stability of the observed performance, we conducted  
 443 a multi-faceted statistical analysis. First, to move beyond simple mean scores,  
 444 we estimated the uncertainty of each model’s performance using bootstrap  
 445 resampling. By creating 10,000 new sets of fold scores via sampling with  
 446 replacement, we constructed 95% confidence intervals (CIs) for the true mean  
 447 macro-Dice score. The results, shown in Table 8 and visualized in Figure 8,  
 448 reveal a pattern of uncertainty. Under the LOOCV protocol, MedSAM’s CI  
 449 ([0.651, 0.734]) is non-overlapping with all other models, providing statistical  
 450 evidence that it is the best-performing architecture. However, among models  
 451 ranked 2nd through 5th, a chain of pairwise CI overlaps prevents confident  
 452 fine-grained ranking: SegFormer ([0.517, 0.633]) overlaps with SAM ([0.466,  
 453 0.527]) and MedSAM+UNet ([0.431, 0.525]), both of which in turn overlap  
 454 with UNet ([0.285, 0.470]). This pattern of overlapping confidence intervals  
 455 indicates that, while MedSAM is statistically separable as the top performer,  
 456 the ranking among the remaining competitive models is not statistically  
 457 robust.  
 458

Table 8: Bootstrap 95% Confidence Intervals for Macro-Dice Score (10 Models)

Model	LOOCV		3-Fold CV	
	Mean	95% CI	Mean	95% CI
<b>MedSAM</b>	<b>0.694</b>	[0.651, 0.734]	<b>0.643</b>	[0.600, 0.688]
SegFormer	0.576	[0.517, 0.633]	0.501	[0.447, 0.554]
SAM	0.496	[0.466, 0.527]	0.425	[0.403, 0.459]
MedSAM+UNet	0.479	[0.431, 0.525]	0.461	[0.395, 0.559]
UNet	0.392	[0.285, 0.470]	0.369	[0.345, 0.416]
UNet-ResNet34	0.373	[0.274, 0.464]	0.268	[0.217, 0.317]
DeepLabV3+	0.304	[0.188, 0.420]	0.335	[0.312, 0.378]
Lightweight UNet	0.249	[0.170, 0.324]	0.213	[0.019, 0.445]
FCN	0.121	[0.080, 0.171]	0.068	[0.054, 0.083]
SegNet	0.105	[0.038, 0.183]	0.098	[0.036, 0.218]

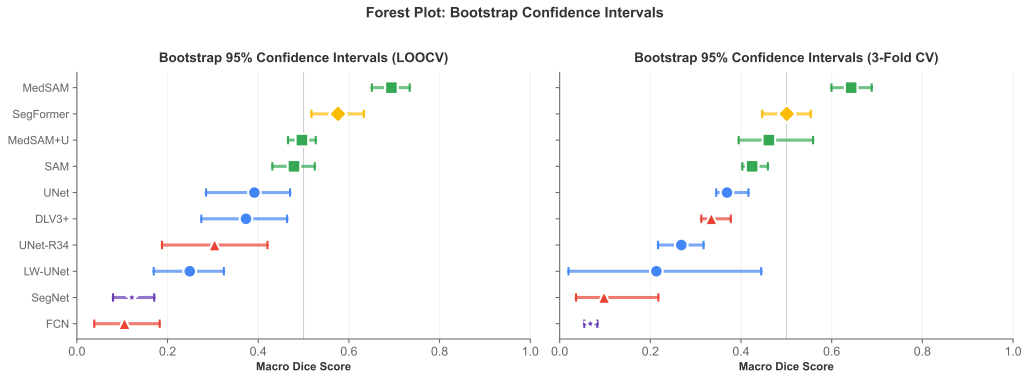


Figure 8: Bootstrap 95% Confidence Intervals for Macro-Dice Score. The plots show the mean Dice score (marker) and the bootstrap 95% confidence interval (10,000 resamples) for each model under 3-Fold CV (left) and LOOCV (right). Under LOOCV, MedSAM’s CI is non-overlapping with all other models, while substantial overlap among 2nd–5th ranked models indicates their fine-grained ranking is not statistically meaningful.

459 Next, we performed a non-parametric rank comparison using the Friedman  
 460 test and Nemenyi post-hoc test, visualized with Critical Difference plots in  
 461 Figure 9. This analysis directly compares the models’ rankings across all folds.  
 462 Both the LOOCV ( $p < 0.001$ ) and 3-Fold CV ( $p = 0.003$ ) analyses yield  
 463 statistically significant Friedman tests, confirming that some rank differences  
 464 exist across the full set of ten models. However, the large critical difference  
 465 threshold in the 3-Fold analysis ( $CD = 7.82$ ) means that the Nemenyi post-  
 466 hoc test connects all top-tier models, providing strong evidence that they  
 467 are statistically indistinguishable from one another despite the overall test  
 468 significance. We note that with only three observations per model, the 3-

469 Fold analysis has limited statistical power for distinguishing ten models; the  
 470 LOOCV analysis—with nine observations per model—provides the more  
 471 reliable basis for pairwise comparison.

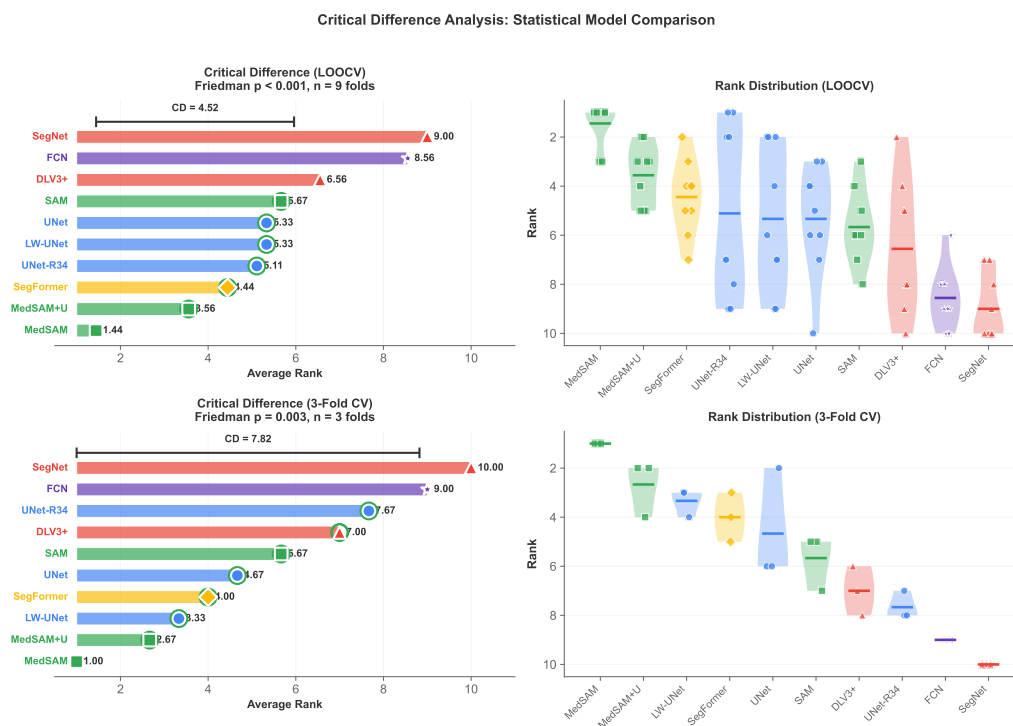


Figure 9: Critical Difference Analysis of Model Ranks. This analysis compares average model ranks under LOOCV (top) and 3-Fold CV (bottom). The plots on the left are Critical Difference (CD) diagrams from the Friedman and Nemenyi tests; models connected by a solid bar are not statistically different. The plots on the right show the distribution of ranks for each model across the validation folds. Both protocols yield significant Friedman tests (LOOCV:  $p < 0.001$ ; 3-Fold CV:  $p = 0.003$ ), confirming overall rank differences. However, the large critical difference in the 3-Fold analysis (CD = 7.82) means the Nemenyi post-hoc test finds no statistically significant pairwise difference among the top-performing models.

472 Finally, to distinguish between statistical significance and practical impor-  
 473 tance, we calculated the effect size (Cohen’s  $d$ ) for all pairwise comparisons  
 474 (Table 9 and Figure 10). This analysis reveals a structured pattern: large  
 475 effect sizes ( $d > 0.8$ ) dominate comparisons between architectural families  
 476 (e.g., foundation models vs. classical CNNs), while within-family comparisons  
 477 (e.g., among CNN variants or among foundation models) yield smaller effects.

478 The rank volatility, detailed in Table 10, confirms that no consistent, repro-  
 479 ducible model hierarchy exists, with high rank standard deviations across  
 480 folds (Table 10). These results indicate that while differences between archi-  
 481 tectural families are meaningful, fine-grained rankings within a family are not  
 482 statistically robust.

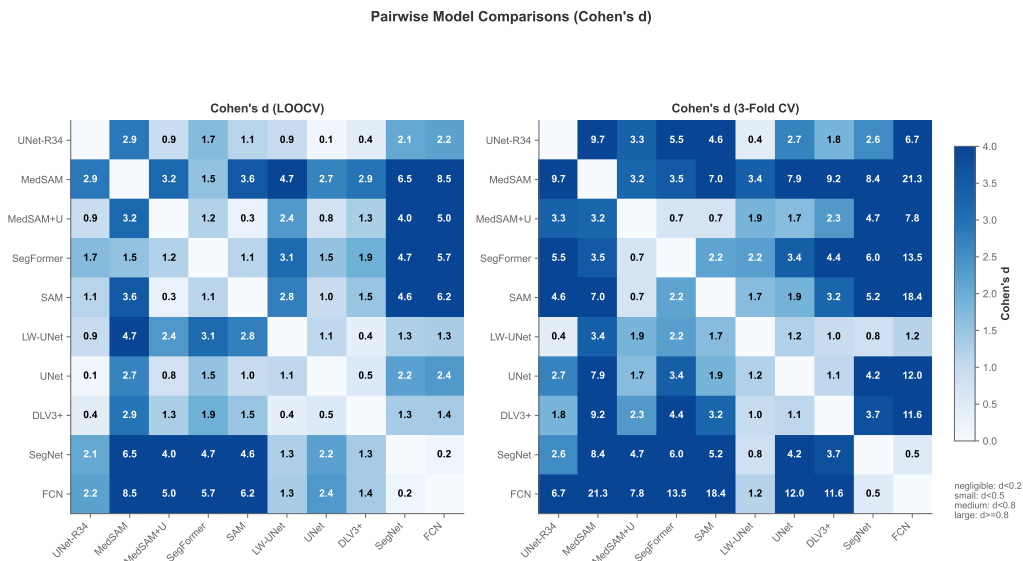


Figure 10: Pairwise Effect Size (Cohen’s d) Matrix for all 10 models. The heatmap visualizes the practical significance of performance differences between model pairs. Large effect sizes ( $d > 0.8$ ) dominate comparisons between architectural families (e.g., foundation models vs. classical architectures), while within-family comparisons yield smaller effects, indicating that fine-grained rankings within a family are not practically meaningful.

### 483 4.3. Qualitative Consistency Despite Quantitative Variability

484 To complement the quantitative analysis, we conducted a multi-modal  
 485 XAI analysis by applying all ten DS1-trained architectures to a representative  
 486 DS2 sample under distribution shift (Figure 11). The 10-row  $\times$  8-column  
 487 grid displays, for each model: the input histological image, the ground truth  
 488 mask, the model’s prediction, an error analysis map, a predictive uncertainty  
 489 map, a morphological gradient, a gradient-based saliency map, and a com-  
 490 bined XAI overlay. The results reveal varying degrees of generalization across  
 491 architectures. UNet and SegFormer produce predictions closest to the ground  
 492 truth, with errors confined primarily to tissue boundaries. MedSAM and

Table 9: Pairwise Effect Size (Cohen’s d) for LOOCV Macro-Dice Scores (10 Models). Positive values indicate the row model outperforms the column model.

vs.	SF	SAM	M-U <sup>b</sup>	UNet	U-R34 <sup>c</sup>	DLV3+ <sup>a</sup>	LW-U <sup>d</sup>	FCN	SegNet
MedSAM	1.55	3.56	3.19	2.69	2.87	2.94	4.72	8.54	6.52
SegFormer	–	1.13	1.21	1.53	1.69	1.95	3.15	5.69	4.70
SAM	–	–	0.29	0.97	1.14	1.49	2.76	6.20	4.57
M-UNet <sup>b</sup>	–	–	–	0.76	0.92	1.30	2.36	5.00	3.99
UNet	–	–	–	–	0.13	0.54	1.08	2.35	2.21
U-R34 <sup>c</sup>	–	–	–	–	–	0.43	0.94	2.20	2.07
DLV3+ <sup>a</sup>	–	–	–	–	–	–	0.36	1.35	1.34
LW-U <sup>d</sup>	–	–	–	–	–	–	–	1.31	1.26
FCN	–	–	–	–	–	–	–	–	0.17

<sup>a</sup> DLV3+: DeepLabV3+    <sup>b</sup> M-UNet: MedSAM+UNet    <sup>c</sup> U-R34: UNet-ResNet34  
<sup>d</sup> LW-U: Lightweight UNet

Table 10: Comparative Rank Instability Analysis Across Cross-Validation Protocols (10 Models)

Model	LOOCV (High Variance)					3-Fold CV (Lower Variance)				
	$\bar{R}$	R	$\sigma$	Top-1	Top-2	$\bar{R}$	R	$\sigma$	Top-1	Top-2
MedSAM	1.00	0.0	0.00	100.0	100.0	1.00	0.0	0.00	100.0	100.0
SegFormer	2.33	1.0	0.50	0.0	66.7	2.33	1.0	0.58	0.0	66.7
SAM	3.78	4.0	1.30	0.0	11.1	4.33	1.0	0.58	0.0	0.0
MedSAM+UNet	4.67	5.0	1.50	0.0	11.1	3.67	4.0	2.08	0.0	33.3
UNet	5.78	7.0	2.11	0.0	0.0	4.67	3.0	1.53	0.0	0.0
UNet-ResNet34	5.78	7.0	2.44	0.0	11.1	7.67	1.0	0.58	0.0	0.0
DeepLabV3+	6.56	6.0	2.24	0.0	0.0	6.00	4.0	2.00	0.0	0.0
LW-UNet	7.44	4.0	1.42	0.0	0.0	6.33	2.0	1.15	0.0	0.0
FCN	8.67	3.0	1.00	0.0	0.0	9.00	0.0	0.00	0.0	0.0
SegNet	9.00	3.0	1.12	0.0	0.0	10.00	0.0	0.00	0.0	0.0

Note:  $\bar{R}$  is the mean rank, R is the rank range (Max - Min),  $\sigma$  is the rank standard deviation, and Top-1/Top-2 show the percentage of folds where the model ranked 1st or within top 2.

493 MedSAM+UNet maintain recognizable segmentation structure with moderate  
494 boundary errors. UNet-ResNet34 and SAM show noisier predictions than  
495 their in-distribution ranking would suggest. DeepLabV3+ exhibits substantial  
496 overprediction of the media class, and LW-UNet displays multi-class confusion.  
497 SegNet and FCN produce scattered, noisy outputs with widespread error and  
498 elevated uncertainty. These results demonstrate that models with comparable  
499 in-distribution performance on DS1 can exhibit substantially different gener-  
500 alization behavior, suggesting that out-of-distribution evaluation provides a  
501 more informative basis for model selection than in-distribution metrics alone.

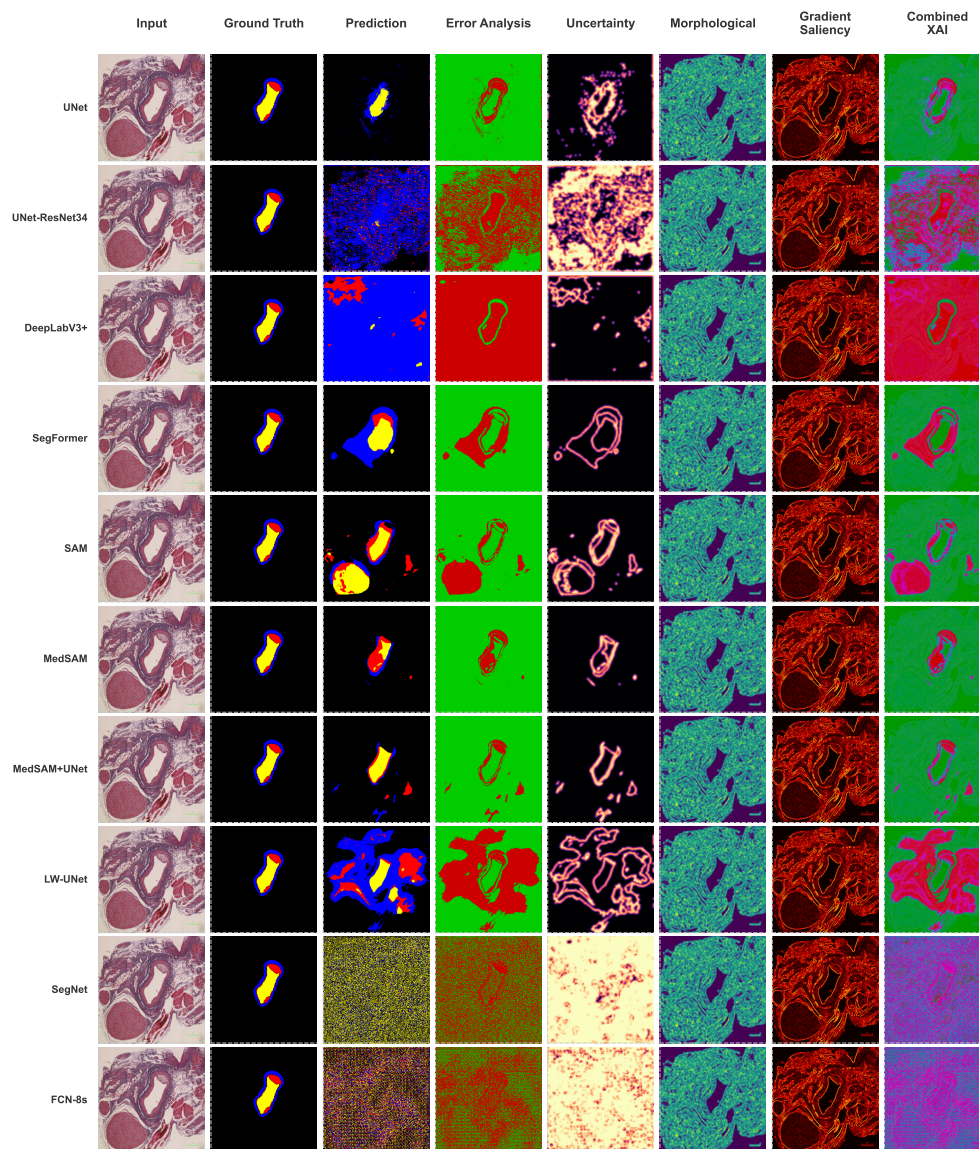


Figure 11: Multi-modal XAI analysis of all ten DS1-trained architectures evaluated on a representative DS2 sample under distribution shift. Each row represents one model; columns show (left to right): input image, ground truth, model prediction, error analysis, predictive uncertainty, morphological gradient, gradient saliency, and combined XAI overlay. Models exhibit varying degrees of generalization, from accurate boundary-level predictions (UNet, SegFormer) to moderate degradation (MedSAM variants) to scattered, noisy outputs (SegNet, FCN)

502 4.4. Quantitative Analysis of XAI Stability

503 To move beyond a purely qualitative assessment of this observation, we  
504 conducted a quantitative analysis of the stability of the models' reasoning.  
505 We used the generated uncertainty maps from the 9 folds of the LOOCV  
506 protocol as a proxy for the model's confidence. For each top-performing  
507 model, we calculated the pixel-wise mean and variance of these maps. The  
508 results for the top-performing models are shown in Figure 12.

509 The mean uncertainty map (Figure 12A) confirms that, on average, the  
510 model's uncertainty is highest at the boundaries between tissue types, which  
511 is clinically expected. Figure 12B displays the variance map, which is nearly  
512 black across the entire image. This indicates a quantitatively minimal variance  
513 in uncertainty across the 9 independent folds. This indicates that while the  
514 quantitative performance metrics were unstable, the model's underlying  
515 reasoning was highly stable and consistent. The instability is therefore  
516 confined to superficial metric noise at the boundaries, not a fundamental  
517 disagreement in model competence.

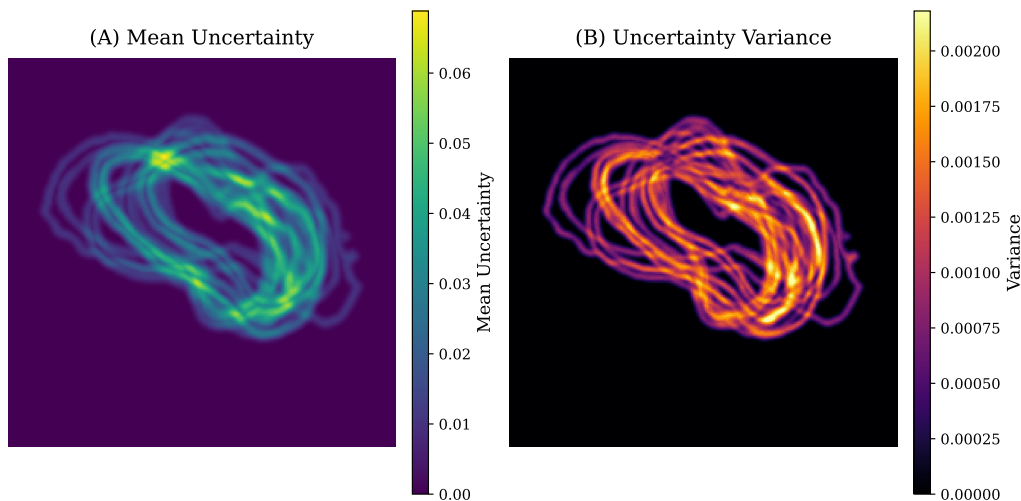


Figure 12: Quantitative XAI Stability Analysis. (A) The mean uncertainty map, averaged across all 9 LOOCV folds, showing that uncertainty is consistently highest at tissue boundaries. (B) The per-pixel variance of uncertainty across the 9 folds. The extremely low variance (dark color) provides quantitative evidence that the model's confidence and reasoning were highly stable, despite fluctuating performance metrics.

518 *4.5. Generalization and Dataset-Dependence of Model Rankings*

519 To test whether the observed benchmark instabilities translate to real-  
520 world deployment scenarios, we conducted two complementary experiments on  
521 an independent generalization dataset (DS2) comprising  $N = 153$  images from  
522 different tissue preparations and magnifications than the training data (DS1).  
523 First, we evaluated all ten DS1-trained models directly on DS2 to assess  
524 robustness under distribution shift (Sections 4.5.1–4.5.3). Second, we trained  
525 all models from scratch on DS2 at varying sample sizes ( $N = 9, 25, 50, 100, 150$ )  
526 and tested on three held-out DS2 images, revealing dataset-specific ranking  
527 hierarchies that differ from those observed on DS1 (Section 4.5.4).

528 Figure 13 illustrates the visual differences between the two datasets. DS1  
529 images exhibit consistent H&E staining with well-defined vessel morphology,  
530 whereas DS2 images display markedly different staining intensity, tissue  
531 architecture, and imaging conditions, confirming that DS2 constitutes a  
532 genuine out-of-distribution challenge rather than a simple extension of DS1.

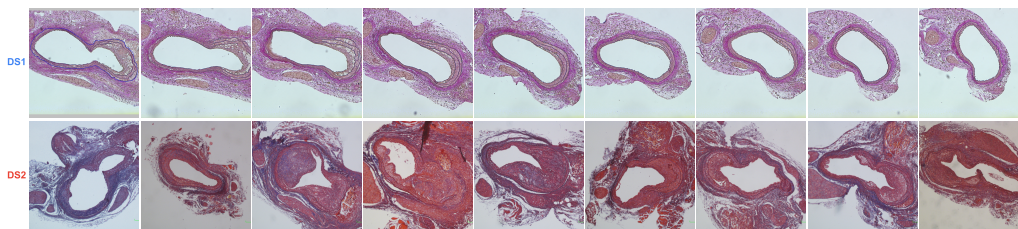


Figure 13: Visual comparison of representative samples from DS1 (top row) and DS2 (bottom row). The two datasets exhibit substantial differences in staining protocols, tissue morphology, and imaging conditions, establishing DS2 as a genuine out-of-distribution evaluation scenario.

533 *4.5.1. Evidence of Distribution Shift*

534 Figure 14 visualizes the distribution shift between DS1 and DS2 using  
535 t-SNE embeddings of features extracted from a pre-trained model. The  
536 clear separation between the two datasets confirms that DS2 represents an  
537 out-of-distribution evaluation scenario.

538 *4.5.2. Sample Size Sensitivity Analysis (Out-of-Distribution)*

539 To understand how model rankings evolve with increasing evaluation  
540 sample sizes under distribution shift, we evaluated DS1-trained models on  
541 nested subsets of DS2 ( $n = 10, 25, 50, 100, 153$  samples). As shown in Figure

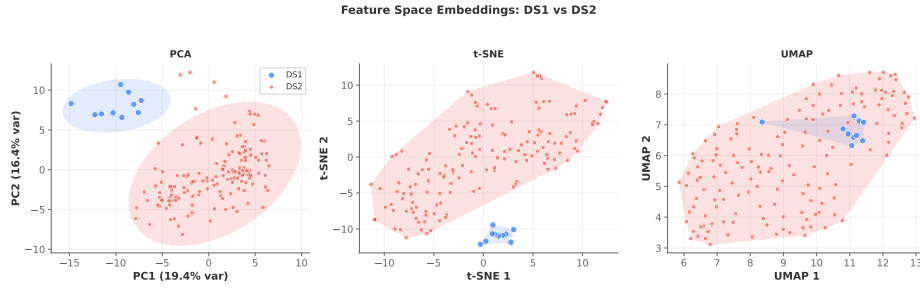


Figure 14: t-SNE visualization of feature embeddings showing distribution shift between training data (DS1,  $N = 9$ ) and generalization data (DS2,  $N = 153$ ). The clear cluster separation indicates substantial domain shift in imaging conditions.

542 15, model rankings are unstable at small sample sizes ( $n \leq 25$ ) but stabilize  
 543 as sample size increases. Foundation models (MedSAM, SAM) maintain  
 544 consistent performance across sample sizes, whereas classical architectures  
 545 (FCN, SegNet) fail to generalize under distribution shift.

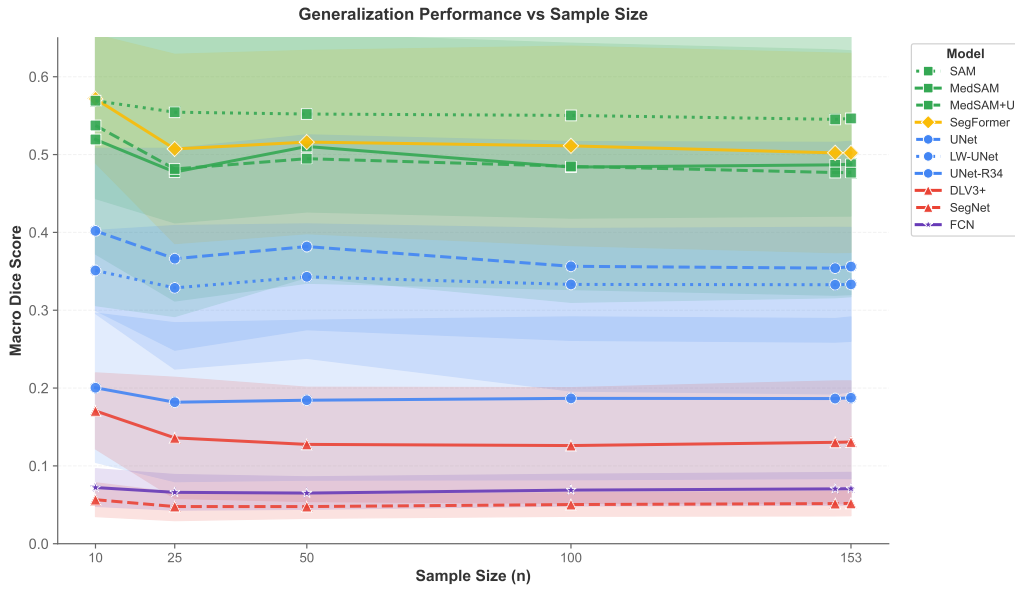


Figure 15: Out-of-distribution generalization performance vs. evaluation sample size. DS1-trained models were evaluated on nested subsets of DS2 ( $n = 10$  to 153). Foundation models maintain stable performance across sample sizes, while classical architectures show persistent failure under distribution shift.

#### 546 4.5.3. Ranking Inversions Under Distribution Shift

547 Model rankings under distribution shift differ from in-distribution rankings.  
548 While MedSAM maintains top performance on both DS1 and DS2, the relative  
549 ranking of other models changes considerably. For instance, SegFormer,  
550 ranked 2nd on DS1, maintains a similar relative position on DS2. However,  
551 UNet variants, which showed competitive DS1 performance, suffer substantial  
552 degradation under shift. Similarly, classical architectures that showed some  
553 learning on DS1 fail almost completely on DS2.

554 Foundation models pre-trained on diverse distributions show better generalization,  
555 suggesting that pre-training strategy may matter more than  
556 architecture choice in data-scarce domains.

557 Figure 16 provides visual evidence of these ranking inversions. The  
558 grid shows DS1-trained models applied to a representative sample of DS2  
559 images spanning the full dataset. SegFormer and DeepLabV3+ maintain  
560 recognizable segmentation structure across samples, whereas SegNet and FCN  
561 produce near-uniform class predictions, confirming their complete failure  
562 under distribution shift. LW-UNet exhibits inconsistent behavior—producing  
563 reasonable segmentations on some samples but fragmented outputs on others—  
564 highlighting the unpredictability of generalization for architectures without  
565 diverse pre-training.

#### 566 4.5.4. In-Distribution Learning Curves on DS2

567 The preceding sections examined how DS1-trained models perform under  
568 distribution shift. A complementary question is: do the same ranking  
569 hierarchies emerge when models are trained *directly* on DS2? To answer  
570 this, we trained all ten architectures on DS2 subsets of increasing size  
571 ( $N = 9, 25, 50, 100, 150$ ), using three held-out DS2 images as a fixed test  
572 set. For  $N \leq 50$ , we used 5-fold cross-validation; for  $N \geq 100$ , an 80/20  
573 train-validation split. All models used the same hyperparameters optimized  
574 on DS1 to isolate dataset effects from tuning effects.

575 Table 11 reports the mean macro Dice score across the three test images  
576 for each model and dataset size. The results show a ranking hierarchy that  
577 differs from DS1.

578 The results demonstrate dataset-specific model rankings. At  $N = 9$  on  
579 DS2, DeepLabV3+ (0.843) and UNet (0.840) lead the ranking—a different  
580 hierarchy from DS1, where MedSAM (0.694) dominated under LOOCV.  
581 Foundation models, which were the clear winners on DS1, rank only 3rd–5th  
582 on DS2 at the same sample size.

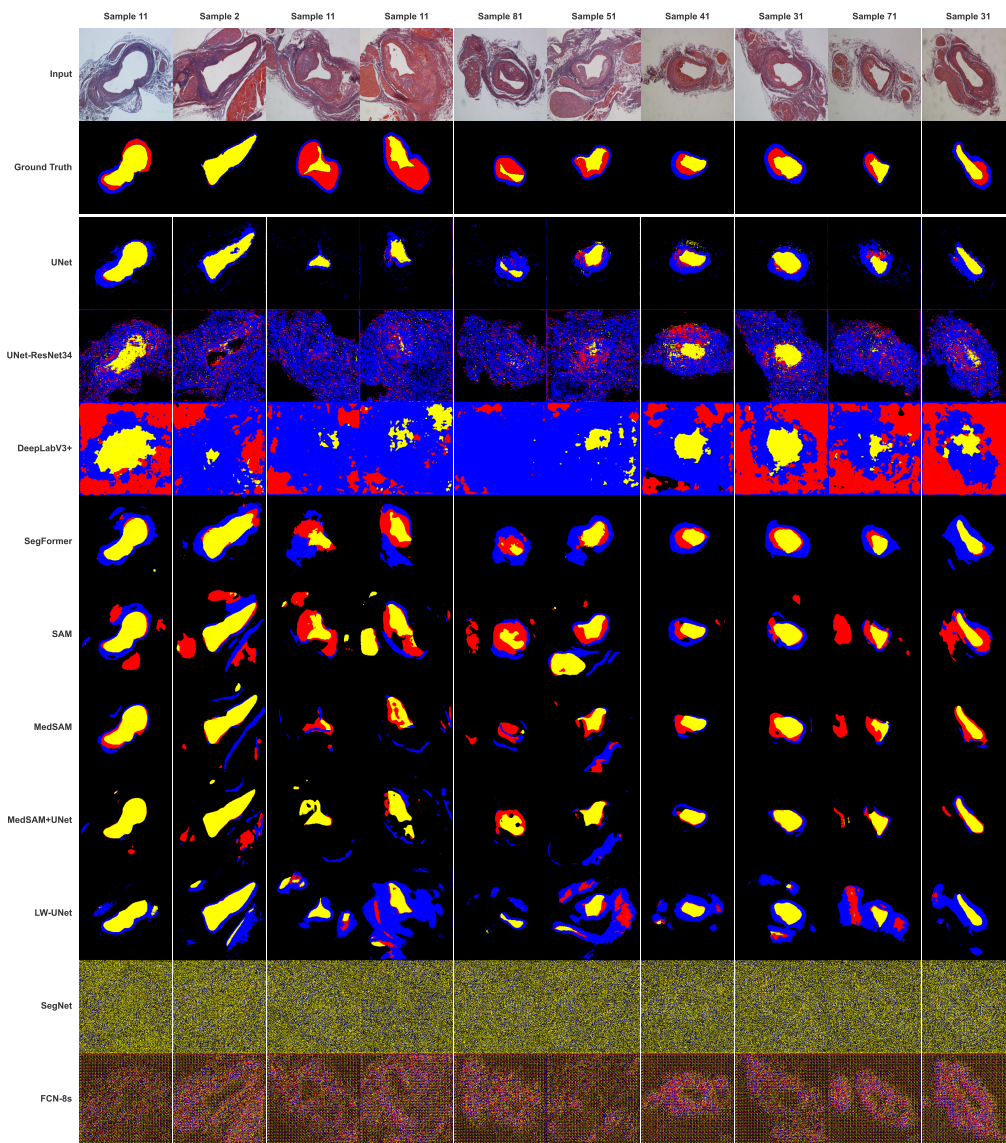


Figure 16: Qualitative inference results of DS1-trained models on DS2 under distribution shift. Rows represent the input image, ground truth, and predictions from each of the ten architectures; columns span representative DS2 samples. Foundation models and top encoders maintain reasonable segmentation, while classical architectures (SegNet, FCN) produce near-uniform or random outputs, visually confirming the ranking inversions observed in quantitative analysis.

Table 11: DS2 In-Distribution Performance: Mean Macro Dice Score. Models trained on DS2 subsets ( $N = 9$  to 150) and tested on three held-out DS2 images. Bold indicates best performance per column.

Model	N=9	N=25	N=50	N=100	N=150
<i>Foundation Models &amp; Transformers</i>					
SegFormer	0.717	0.657	0.847	0.896	<b>0.906</b>
MedSAM	0.779	0.839	0.846	0.857	0.870
MedSAM+UNet	0.812	0.833	0.839	0.859	0.864
SAM	0.821	0.792	0.836	0.859	0.866
<i>Modern CNNs</i>					
UNet	0.840	0.820	0.837	0.889	0.886
UNet-ResNet34	0.769	<b>0.862</b>	<b>0.865</b>	0.876	0.892
DeepLabV3+	<b>0.843</b>	0.838	0.852	<b>0.891</b>	0.896
Lightweight UNet	0.605	0.810	0.787	0.879	0.904
<i>Classical Architectures</i>					
FCN	0.450	0.549	0.562	0.726	0.828
SegNet	0.051	0.067	0.618	0.634	0.586

583 The scaling behavior of different architectures further distinguishes them.  
584 SegFormer exhibits classic data-hungry transformer behavior: it ranks near  
585 the bottom at  $N = 9$  (0.717, rank 7/10) but rises to the top at  $N = 150$  (0.906,  
586 rank 1/10), a +26% improvement. Similarly, Lightweight U-Net improves  
587 from 0.605 to 0.904 (+49%). This contrasts with foundation models, which  
588 show only modest improvement (+5–12% from  $N = 9$  to  $N = 150$ ). SAM,  
589 MedSAM, and MedSAM+UNet converge to a narrow band around 0.86–0.87  
590 at  $N = 150$ , despite starting from different points. Their pre-training provides  
591 a strong initialization that helps at small  $N$  but offers diminishing returns as  
592 in-distribution data grows.

593 At the other end of the spectrum, classical architectures perform poorly.  
594 SegNet achieves only 0.051 Dice at  $N = 9$ , approaching chance-level performance—  
595 and never exceeds 0.634 even at  $N = 100$ . FCN similarly struggles (0.450 at  
596  $N = 9$ ), confirming that these older architectures lack the inductive biases  
597 necessary for learning from minimal data. Finally, the data also independently  
598 confirm the stability threshold identified in the DS1 analysis (Section 4.6):  
599 rankings are volatile at  $N \leq 25$  (e.g., SegFormer drops from rank 7 to rank 8  
600 between  $N = 9$  and  $N = 25$ ) but stabilize at  $N \geq 50$ .

601 Figure 17 provides a six-panel analysis of these scaling dynamics. The  
602 rank trajectory (panel a) and rank heatmap (panel b) visualize the frequent  
603 crossing patterns at small  $N$  that gradually converge as data increases. The

604 performance-vs-dataset-size small multiples (panel c) reveal architecture-  
605 specific scaling curves, with SegFormer and LW-UNet showing steep improve-  
606 ment while foundation models plateau early. The confidence interval width  
607 (panel d) decreases monotonically with  $N$ , confirming that larger samples  
608 yield tighter estimates. Kendall’s  $\tau$  correlation with the final ranking (panel  
609 e) reaches near-perfect agreement by  $N = 100$ , and the number of rank  
610 swaps between consecutive dataset sizes (panel f) drops sharply after  $N = 50$ ,  
611 independently confirming the stability threshold identified on DS1. The quali-  
612 tative predictions in Figure 18 provide visual confirmation: at  $N = 9$ , several  
613 models produce fragmented or missing segmentations, while at  $N = 150$  most  
614 models achieve visually accurate delineation of all tissue classes.

#### 615 4.5.5. *Synthesis: Domain vs. Architecture*

616 The preceding experiments reveal that model selection cannot be captured  
617 by a single leaderboard. The out-of-distribution results (Sections 4.5.1–4.5.3)  
618 demonstrate that foundation models perform better when deployed on unseen  
619 data distributions, maintaining reasonable Dice scores where conventional  
620 architectures degrade. However, the in-distribution DS2 results (Section  
621 4.5.4) show that this advantage is specifically about pre-training diversity, not  
622 inherent architectural superiority. When trained directly on DS2, conventional  
623 architectures like DeepLabV3+ and SegFormer match or exceed foundation  
624 model performance.

625 This distinction has direct practical implications. When distribution shift  
626 is expected, foundation models offer a clear advantage through their diverse  
627 pre-training. Distribution shifts often occur in clinical deployment where  
628 training data may come from one institution and test data from another. When  
629 sufficient in-distribution data is available ( $N \geq 50$ ), task-specific architectures  
630 like SegFormer or DeepLabV3+ can achieve superior performance with far  
631 fewer computational resources. The ‘best’ model is therefore not a fixed  
632 property of an architecture but a function of both the available data and the  
633 expected deployment conditions.

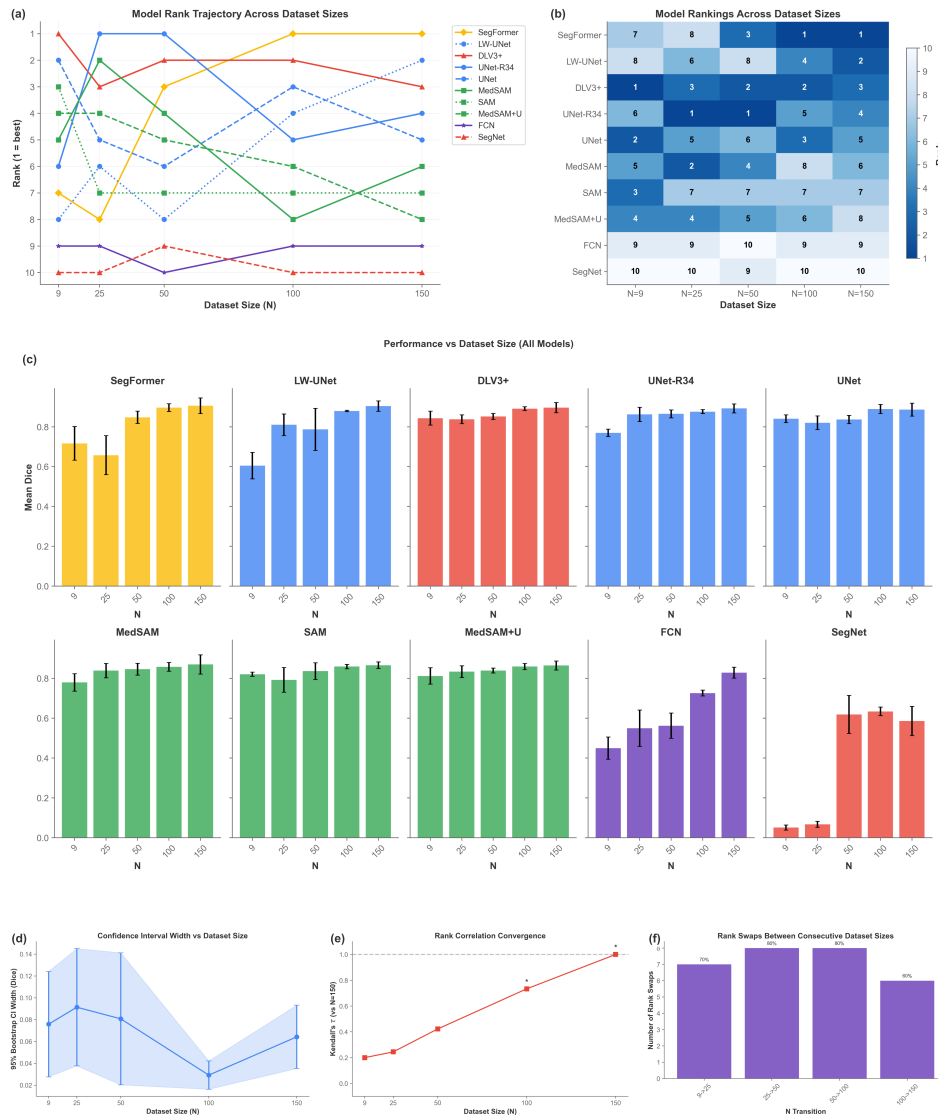


Figure 17: Comprehensive DS2 in-distribution analysis across dataset sizes ( $N = 9$  to  $150$ ). (a) Rank trajectory showing frequent rank crossings at small  $N$  that stabilize at  $N \geq 50$ . (b) Rank heatmap providing a compact view of all models' rankings across dataset sizes. (c) Performance vs. dataset size for each architecture, revealing distinct scaling behaviors. (d) Bootstrap 95% CI width decreasing with sample size, confirming tighter estimates at larger  $N$ . (e) Kendall's  $\tau$  rank correlation convergence toward the final ( $N = 150$ ) ranking. (f) Number of rank swaps between consecutive dataset sizes, with a sharp decline after  $N = 50$  confirming the stability threshold.

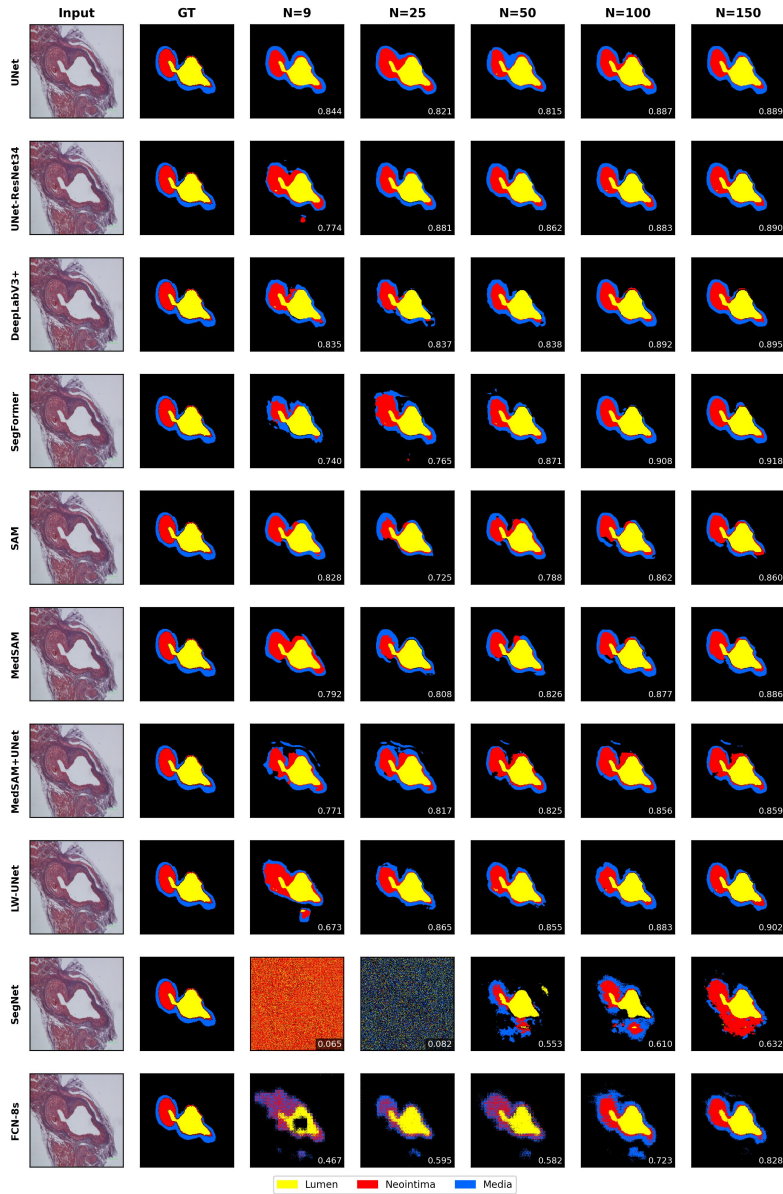


Figure 18: Qualitative predictions on a held-out DS2 test image across all ten models and dataset sizes ( $N = 9$  to 150). Each row shows one architecture with Dice scores annotated per panel. Columns represent increasing training data. At  $N = 9$ , several models produce fragmented segmentations; by  $N = 150$ , most models converge to visually accurate predictions. Classical architectures (SegNet, FCN) show persistent failure even at larger  $N$ , while SegFormer exhibits the largest improvement from 0.740 to 0.918.

634 *4.6. Ablation Studies*

635 To systematically quantify sources of experimental variance beyond data  
 636 splits, we conducted 190+ ablation experiments examining augmentation  
 637 strategies, input resolution, random seed effects, and hyperparameter sensi-  
 638 tivity.

639 *4.6.1. Data Augmentation and Seed Stability*

640 Figure 19 summarizes two important sources of variance. First, we tested  
 641 10 different augmentation presets (10 experiments  $\times$  10 models = 100 total  
 642 runs), ranging from no augmentation to aggressive geometric and color trans-  
 643 formations. The variance introduced by augmentation choice was substantial—  
 644 for some models, the difference between best and worst augmentation preset  
 645 exceeded 0.15 Dice points. Second, we trained each model with 5 different ran-  
 646 dom seeds (5 seeds  $\times$  10 models = 50 runs). Seed-induced variance was smaller  
 647 than augmentation variance but still meaningful, with standard deviations  
 648 ranging from 0.03 to 0.15 Dice points across models, with foundation models  
 649 exhibiting the least seed sensitivity (0.03–0.06) and classical/transformer  
 650 architectures the most (up to 0.11 for SegNet and 0.15 for SegFormer).

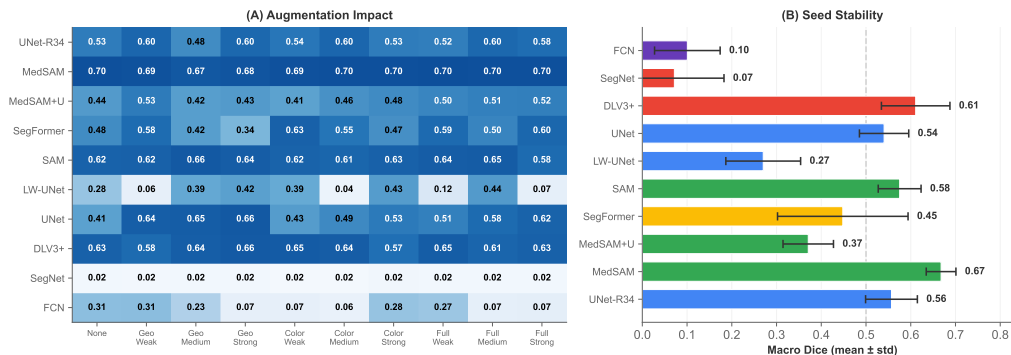


Figure 19: Ablation study: Data augmentation and seed stability. Left: Performance variance across 10 augmentation presets per model. Right: Performance variance across 5 random seeds per model. Both sources introduce variance comparable to architectural differences.

651 Beyond demonstrating that augmentation introduces variance, the heatmap  
 652 in Figure 19A reveals architecture-specific augmentation preferences that yield  
 653 practical guidance for practitioners.

654 Foundation models are augmentation-robust: MedSAM is nearly insen-  
 655 sitive to augmentation choice (Dice range: 0.67–0.70 across all 10 presets),

656 achieving competitive performance even without augmentation. This robust-  
657 ness likely stems from the diverse pre-training distribution of the ViT encoder,  
658 which already encodes invariances that augmentation would otherwise need  
659 to teach. By contrast, CNNs benefit most from geometric augmentation.  
660 UNet improved from 0.41 (no augmentation) to 0.66 with strong geometric  
661 transforms (+61%), and DeepLabV3+ peaked at 0.66 with the same pre-  
662 set. Geometric augmentations (rotations, flips, affine transforms) effectively  
663 expand the limited training distribution for architectures that lack built-in  
664 spatial invariances; for CNN-based architectures, we recommend geometric  
665 augmentation at medium-to-strong intensity.

666 SegFormer uniquely prefers color augmentation. Unlike CNNs, it achieved  
667 its best performance with weak color augmentation (0.63) while strong geo-  
668 metric transforms degraded it to 0.34. The hierarchical transformer encoder,  
669 which already captures spatial relationships through self-attention, benefits  
670 more from staining variability simulation than from geometric perturbations;  
671 we therefore recommend color-based augmentation for transformer architec-  
672 tures. Finally, lightweight architectures are highly sensitive to augmentation  
673 choice: Lightweight U-Net exhibited substantial instability, with Dice scores  
674 ranging from 0.04 (color medium) to 0.44 (full medium)—a  $10\times$  performance  
675 ratio from augmentation choice alone. This suggests that under-parameterized  
676 models lack the capacity to simultaneously learn the task and cope with aggres-  
677 sive augmentation, making augmentation selection important for lightweight  
678 deployment scenarios.

679 These findings provide architecture-family-specific augmentation recom-  
680 mendations rather than treating augmentation as a universal preprocessing  
681 step.

#### 682 *4.6.2. Resolution Sensitivity*

683 We evaluated all models at four input resolutions (128, 256, 512, 1024  
684 pixels), with results shown in Figure 20. Foundation models generally benefit  
685 from higher resolutions (512–1024px), with SAM peaking at 1024px and  
686 MedSAM at 512px. CNN architectures show mixed resolution preferences:  
687 some (UNet, DeepLabV3+) peak at 1024px, while others (UNet-ResNet34,  
688 Lightweight UNet) perform best at 512px. SegFormer uniquely peaks at  
689 128px, suggesting its transformer architecture captures sufficient context even  
690 at low resolution. Ultra-low resolution (128px) reduces performance for most  
691 models.

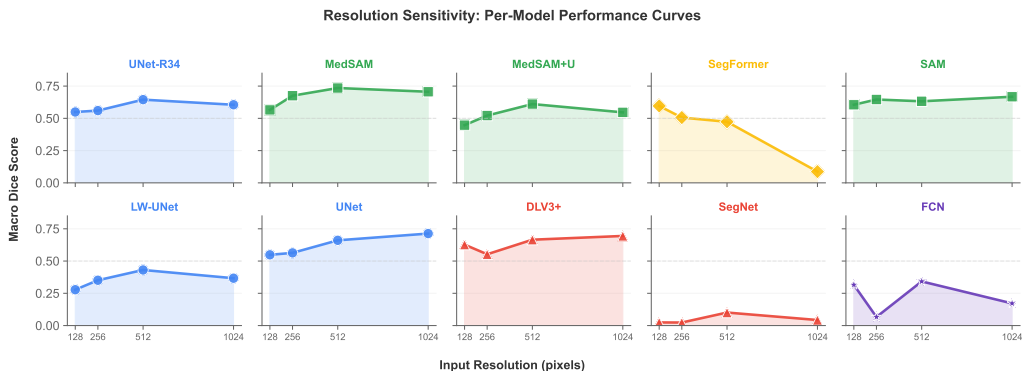


Figure 20: Resolution sensitivity analysis. Performance across input resolutions (128–1024 pixels) reveals architecture-specific optimal operating points, with no universal resolution preference across model families.

#### 692 4.6.3. Hyperparameter Sensitivity

693 To quantify how hyperparameter choices influence model performance,  
 694 we analyzed 100-trial optimization sweeps for each architecture, examining  
 695 learning rate, batch size, optimizer, scheduler, weight decay, and dropout.  
 696 Figures 21, 22, and 23 present this analysis, revealing architecture-specific  
 697 sensitivities that further complicate the notion of fair model comparison.

698 Three findings emerge from this analysis. First, models exhibit sub-  
 699 stantially different optimal learning rate ranges. Foundation models (SAM,  
 700 MedSAM) show narrow optimal ranges around  $10^{-4}$ , while classical architec-  
 701 tures tolerate broader ranges, suggesting that under-tuned learning rates may  
 702 systematically disadvantage certain architectures. Second, relative hyperpa-  
 703 rameter importance varies across architectures. Learning rate dominates for  
 704 most models, but foundation models show greater sensitivity to weight decay  
 705 (affecting fine-tuning behavior), while classical architectures are more sensi-  
 706 tive to dropout and optimizer choice. Third, AdamW generally outperforms  
 707 alternatives for transformer-based models (SegFormer, foundation models),  
 708 while classical architectures show less optimizer sensitivity. This last finding  
 709 suggests that optimizer choice, rarely reported in comparative studies, may  
 710 introduce systematic bias.

711 Figure 22 demonstrates that the relative importance of each hyperparam-  
 712 eter also varies by architecture beyond optimal ranges. Moreover, the choice  
 713 of optimizer introduces an additional, often unreported source of systematic  
 714 bias as illustrated in Figure 23.

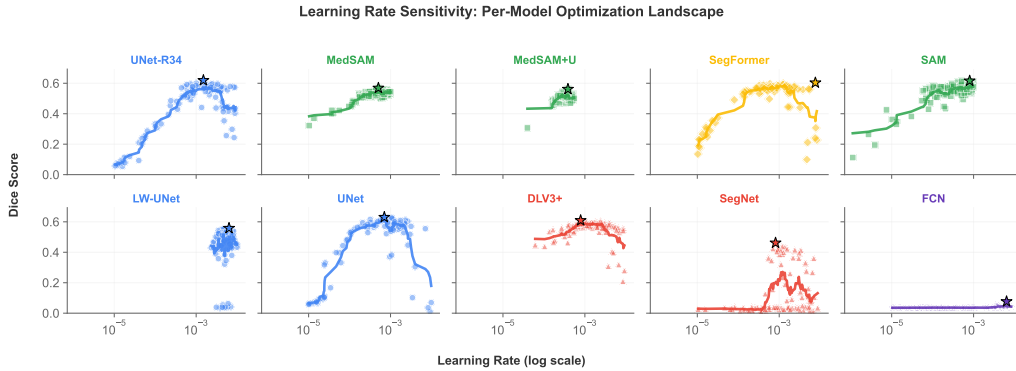


Figure 21: Learning Rate Sensitivity Analysis. Performance distributions across learning rate ranges for each model, revealing architecture-specific optimal ranges. Foundation architectures (SAM, MedSAM) show narrow optimal ranges around  $10^{-4}$ , while classical architectures tolerate broader ranges.

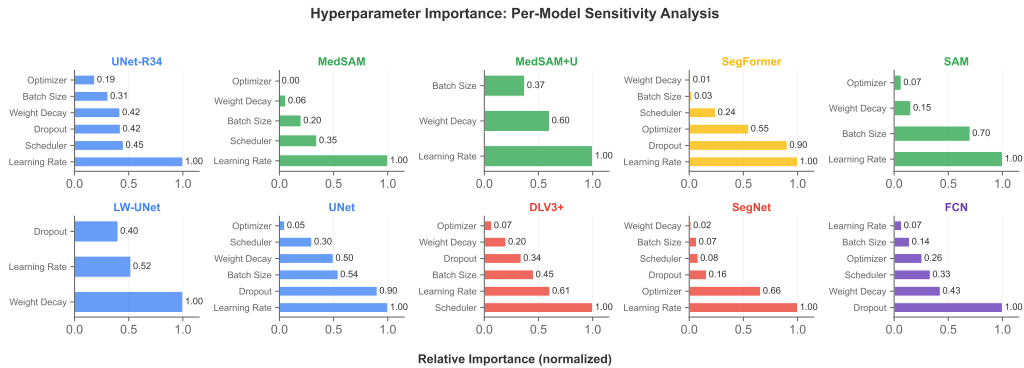


Figure 22: Hyperparameter Importance Analysis. Relative importance of hyperparameters computed via Spearman correlation with validation performance. Learning rate dominates for most models, but foundation models show greater sensitivity to weight decay, while classical architectures are more sensitive to dropout and optimizer choice.

#### 715 4.6.4. Computational Efficiency

716 Beyond segmentation accuracy, practical deployment requires considering  
 717 computational costs. Table 12 summarizes the efficiency metrics for all  
 718 architectures, showing a wide range of resource requirements.

719 The efficiency analysis reveals that SegFormer, with only 3.7M parameters,  
 720 achieves competitive performance with  $25\times$  fewer parameters than founda-  
 721 tion models (Figure 24). This makes it suitable for resource-constrained

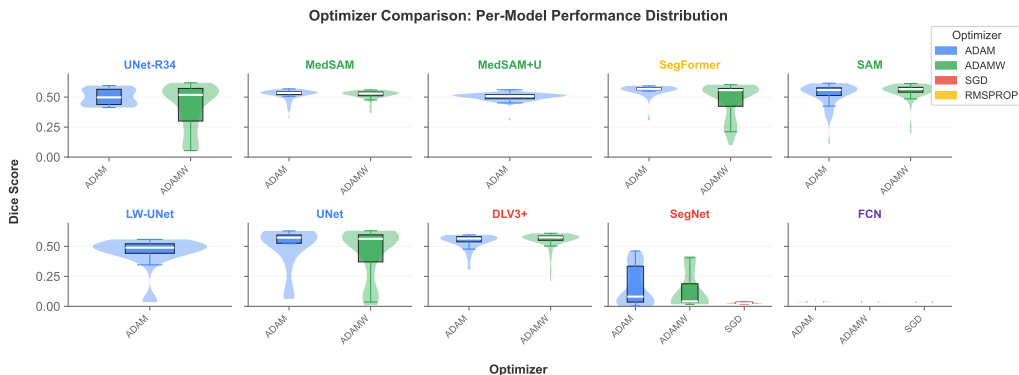


Figure 23: Optimizer Comparison. Performance distributions across Adam, AdamW, SGD, and RMSprop for each architecture. AdamW generally outperforms alternatives for transformer-based models, while classical architectures show less optimizer sensitivity.

Table 12: Computational Efficiency Comparison. Metrics measured on NVIDIA V100 GPU with  $256 \times 256$  input resolution. Training time is the mean wall-clock time per LOOCV fold.

Model	Params (M) ↓	FLOPs (T) ↓	VRAM (GB) ↓	Latency (ms) ↓	FPS ↑	Train (s) ↓	LOOCV Dice ↑
SegFormer	<b>3.7</b>	<b>0.49</b>	<b>0.07</b>	9.3	108	99	0.576
Lightweight UNet	7.8	1.02	0.08	<b>2.7</b>	<b>365</b>	86	0.249
UNet-ResNet34	24.4	3.20	0.13	5.7	175	104	0.373
DeepLabV3+	26.7	3.50	0.14	7.6	132	96	0.304
SegNet	29.5	3.86	0.20	5.7	174	<b>71</b>	0.105
UNet	31.4	4.11	0.25	7.7	131	145	0.392
MedSAM	95.8	12.56	2.71	121.8	8.2	259	<b>0.694</b>
MedSAM+UNet	95.7	12.55	2.71	137.3	7.3	173	0.479
SAM	97.3	12.75	2.72	120.9	8.3	311	0.496
FCN	134.4	17.61	5.40	14.3	70	90	0.121

722 deployment, despite slightly lower absolute performance.

#### 723 4.6.5. Boundary-Sensitive Metrics

724 Beyond overlap-based metrics (Dice, IoU), accurate boundary delineation  
725 is important for histological analysis where tissue layer boundaries inform  
726 diagnosis. We therefore evaluated all models using boundary-sensitive metrics  
727 on the generalization dataset (DS2,  $N = 153$ ): the 95th percentile Hausdorff  
728 Distance (HD95) and Average Symmetric Surface Distance (ASSD) [66, 67].  
729 We focus on DS2 boundary metrics as they reflect the more clinically relevant  
730 scenario of performance under distribution shift, where boundary accuracy  
731 differences between architectures are most pronounced.

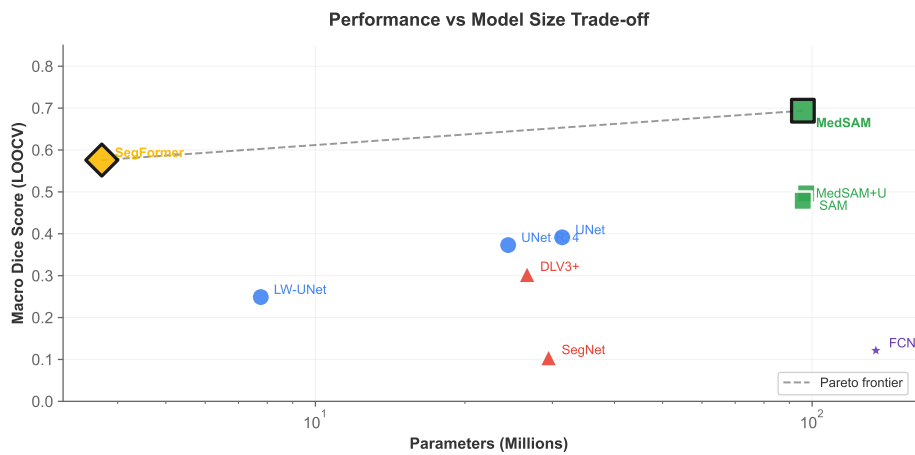


Figure 24: Performance-Efficiency Trade-off. Bubble size represents model parameter count. SegFormer achieves the best balance of high performance with low computational cost, while foundation models (SAM, MedSAM) incur significant overhead.

732 Table 13 presents the boundary metric results. Lower values indicate better  
 733 boundary accuracy. The results show a hierarchy: foundation models and  
 734 transformers achieve the lowest boundary errors, with SegFormer achieving  
 735 HD95 of 27.2 pixels and ASSD of 8.3 pixels. Classical architectures (FCN,  
 736 SegNet) exhibit substantially larger boundary errors (HD95 > 110 pixels),  
 737 confirming their failure to learn precise tissue boundaries under distribution  
 738 shift. Notably, the Dice scores under distribution shift are substantially lower  
 739 than the in-distribution DS1 LOOCV values (e.g., SAM achieves 0.546 on  
 740 DS2 vs. 0.496 on DS1, while DeepLabV3+ drops from 0.304 to 0.131), and  
 741 the ranking order changes—SAM leads in generalization Dice despite ranking  
 742 3rd on DS1.

743 SegFormer achieves the best boundary metrics (HD95, ASSD) despite  
 744 not having the highest generalization Dice score, suggesting that lightweight  
 745 transformer architectures may learn more precise boundary representations  
 746 than larger foundation models. The ranking inversion between boundary  
 747 and overlap metrics—SAM leads in Dice while SegFormer leads in boundary  
 748 precision—reinforces our recommendation to consider multiple evaluation  
 749 dimensions when selecting models for clinical deployment.

Table 13: Boundary-Sensitive Metrics on Generalization Dataset (DS2,  $N = 153$ ). HD95 = 95th percentile Hausdorff Distance; ASSD = Average Symmetric Surface Distance. Lower values indicate better boundary accuracy. Values are mean  $\pm$  std across foreground classes (Lumen, Neointima, Media).

Model	HD95 (px) $\downarrow$	ASSD (px) $\downarrow$	Dice $\uparrow$	Boundary Rank
<i>Foundation Models &amp; Transformers</i>				
<b>SegFormer</b>	<b>27.2 <math>\pm</math> 18.4</b>	<b>8.3 <math>\pm</math> 6.6</b>	0.502	1
SAM	47.9 $\pm$ 25.1	12.5 $\pm$ 7.7	<b>0.546</b>	2
MedSAM	33.1 $\pm$ 21.2	9.9 $\pm$ 7.8	0.488	3
MedSAM+UNet	38.9 $\pm$ 22.0	10.9 $\pm$ 7.6	0.477	4
<i>Modern CNNs</i>				
UNet	56.4 $\pm$ 24.9	16.8 $\pm$ 18.1	0.356	5
Lightweight UNet	56.7 $\pm$ 16.6	19.9 $\pm$ 6.8	0.333	6
UNet-ResNet34	88.7 $\pm$ 11.4	32.6 $\pm$ 5.7	0.188	7
DeepLabV3+	106.9 $\pm$ 13.8	54.0 $\pm$ 12.5	0.131	8
<i>Classical Architectures</i>				
SegNet	111.7 $\pm$ 8.8	50.1 $\pm$ 7.4	0.052	9
FCN	113.4 $\pm$ 8.0	53.0 $\pm$ 6.4	0.071	10

## 750 5. Discussion

751 In this study, we comparatively evaluated ten different deep learning  
 752 segmentation models on a small cardiovascular histology dataset to assess  
 753 carotid artery segmentation for benchmarking. Our results based on extensive  
 754 hyperparameter optimization and ablation studies showed that the concept  
 755 of a single best model is unreliable and observed model rankings are largely  
 756 influenced by evaluation protocols, statistical variability, and experimental  
 757 settings.

### 758 5.1. Benchmark Reliability in Low-Data Settings

759 The central finding of our work is that benchmark leadership is an artifact  
 760 of the evaluation protocol, not an intrinsic property of a model. This insta-  
 761 bility is driven by a combination of the well-known bias-variance trade-off  
 762 in validation strategies and the brittleness of hyperparameter optimization.  
 763 As shown in our results in Figure 7, the high variance of LOOCV created  
 764 unstable model rankings sensitive to single data points, while the higher bias  
 765 of 3-Fold CV may have unfairly penalized more complex models.

766 This phenomenon provides strong empirical evidence for the illusion of  
 767 control cognitive bias in machine learning research [20]. The substantial  
 768 effort invested in rigorous hyperparameter tuning creates overconfidence

769 that randomness has been controlled and a definitive outcome has been  
770 reached. However, our findings show that even with optimal configurations,  
771 performance remains sensitive to the particular composition of data splits.  
772 This highlights a limitation in the standard practice of declaring a state-  
773 of-the-art model based on a single leaderboard, particularly when data is  
774 scarce.

### 775 *5.2. Limitations of Standard Evaluation Metrics*

776 On in-distribution data, quantitative metrics suggest variability among  
777 top-performing models. However, our multi-modal XAI analysis under dis-  
778 tribution shift as shown in Figure 11 reveals that these differences become  
779 more pronounced on unseen data: models that appeared statistically com-  
780 parable on DS1 exhibit varying degrees of generalization on DS2, ranging  
781 from accurate predictions to substantial degradation. This suggests that  
782 in-distribution metric variability likely stems from minor, pixel-level boundary  
783 disagreements that may not be clinically meaningful [34], while generalization  
784 capacity represents a more informative criterion for distinguishing between  
785 architectures.

### 786 *5.3. Generalization of Model Selection*

787 Our generalization experiments on DS2 (Section 4.5) provide the most prac-  
788 tically relevant insights. The ranking inversions observed under distribution  
789 shift—where models that performed well on DS1 failed on DS2—demonstrate  
790 that in-distribution benchmarking may not predict real-world deployment  
791 success. Foundation models (MedSAM, SAM) maintained reasonable perfor-  
792 mance under distribution shift, while classical architectures (FCN, SegNet)  
793 showed severe degradation.

794 Crucially, the in-distribution DS2 experiment (Section 4.5.4) demonstrates  
795 that these ranking inversions are not merely an artifact of distribution shift.  
796 When models are trained directly on DS2, the ranking hierarchy differs  
797 completely from DS1—DeepLabV3+ and UNet lead at  $N = 9$  on DS2,  
798 whereas MedSAM led on DS1. Foundation models lose their advantage  
799 when in-distribution training data is available, confirming that rankings are  
800 dataset-specific, not just protocol-specific.

801 A methodological caveat is that our DS2 experiments reused hyperpa-  
802 rameters optimized on DS1, which could systematically favor architectures  
803 whose optimal configurations transfer well across datasets. This design choice  
804 reflects a realistic deployment scenario, but the observed DS2 rankings may

805 partly reflect hyperparameter transferability rather than pure architectural  
806 capacity.

807 This finding has direct implications for model selection in medical imaging.  
808 Rather than optimizing for marginal gains on a single benchmark, practitioners  
809 should prioritize robustness under distribution shift, where foundation models  
810 with diverse pre-training show superior generalization. Models should also  
811 be evaluated across varying sample sizes to assess stability. Furthermore,  
812 independent test sets from different institutions or imaging protocols are  
813 essential for domain-specific validation. Finally, rankings established on one  
814 dataset should not be assumed to transfer to another, even at matching  
815 sample sizes, underscoring the need for cross-dataset validation.

#### 816 5.4. *The Statistical Stability of Performance Rankings*

817 This study was conducted in a very low-data environment ( $N = 9$ ), which  
818 serves as a challenging test case for benchmarking practices. Our sample size  
819 sensitivity analysis in Figure 17a–b empirically maps the phase transition  
820 from instability to stability, with the rank correlation convergence in Figure  
821 17e and rank swap analysis in Figure 17f confirming that rankings become  
822 reproducible only at  $N \geq 50$  samples.

823 The volatile model rankings at small sample sizes, particularly under  
824 the high-variance LOOCV protocol as shown in Figure 7, are an expected  
825 consequence of high variance estimators. However, our ablation studies  
826 (Section 4.6) reveal that variance arises from multiple sources—not just data  
827 splits, but also augmentation choices, random seed, and hyperparameter  
828 selection. Even with controlled randomness, the small dataset size introduces  
829 substantial variance.

830 Importantly, the DS2 in-distribution experiment (Section 4.5.4) indepen-  
831 dently confirms this stability threshold: rankings on DS2 are volatile at  
832  $N \leq 25$  but stabilize at  $N \geq 50$ , providing cross-dataset validation of this  
833 finding.

834 We therefore propose evidence-based guidelines: confident model selection  
835 requires a minimum of 50–100 samples for stable rankings, and performance  
836 claims should always be accompanied by bootstrap confidence intervals and  
837 cross-protocol validation. This guidance is essential for researchers in rare  
838 diseases or other data-limited domains.

839 *5.5. Understanding Hyperparameter Sensitivity Across Architectures*

840 Our hyperparameter optimization experiments revealed architecture-specific  
841 sensitivities that merit mechanistic explanation, as these insights are crucial  
842 for practitioners seeking to replicate our findings or adapt these models to  
843 new domains.

844 **Learning Rate Preferences:** The optimal learning rate varies substantially  
845 across architecture families. Transformer-based models (SegFormer, MiT-B0  
846 encoder) converged optimally with learning rates around  $6 \times 10^{-5}$ , consistent  
847 with established fine-tuning guidelines for pre-trained vision transformers  
848 [68]. This conservative learning rate preserves the rich representations learned  
849 during pre-training while allowing task-specific adaptation. In contrast, foun-  
850 dation models (SAM, MedSAM) using parameter-efficient fine-tuning (Norm  
851 Tuning) favored slightly higher learning rates ( $1 \times 10^{-4}$  to  $3 \times 10^{-4}$ ), as the  
852 limited number of trainable normalization parameters require sufficient gradi-  
853 ent magnitude to learn meaningful domain adaptations without disrupting  
854 the frozen backbone [69]. Classical CNNs trained from scratch exhibited  
855 broader tolerance, performing reasonably across a wider learning rate range.

856 **Loss Function Selection:** The dominance of Focal-Dice loss across all archi-  
857 tectures (Table 2) is explained by its dual mechanism for handling the class  
858 imbalance inherent in our dataset, where tissue layers occupy substantially  
859 different pixel proportions. The Dice component provides spatial overlap  
860 sensitivity that is inherently robust to class imbalance, while the Focal com-  
861 ponent down-weights easy examples to focus learning on difficult boundary  
862 pixels [59]. This combination addresses both class imbalance (different tissue  
863 proportions) and difficulty imbalance (easy interior vs. hard boundary pixels),  
864 making it particularly effective for medical image segmentation where accurate  
865 boundary delineation is clinically critical.

866 **Optimizer Choice:** AdamW consistently outperformed alternatives for  
867 transformer-based architectures, likely due to its decoupled weight decay  
868 regularization which prevents overfitting in over-parameterized models [70].  
869 Classical CNNs showed less optimizer sensitivity, with Adam and SGD with  
870 momentum producing comparable results. These findings underscore that  
871 optimizer choice can systematically bias results and should be tuned indepen-  
872 dently for each architecture family.

873 **6. Conclusion**

874 In this study, we conducted a comprehensive evaluation of ten deep  
875 learning segmentation models on a realistic, data-scarce medical dataset of  
876 cardiovascular histology images, to assess their performance in segmenting  
877 carotid artery structures. Based on the evaluations through extensive hy-  
878 perparameter optimization runs and ablation experiments, we empirically  
879 observed that model rankings depend on the evaluation protocol, statistical  
880 noise, and experimental conditions rather than reflecting true architectural  
881 advantage.

882 Our findings indicate that bootstrap confidence intervals for top-performing  
883 models substantially overlap, indicating statistical indistinguishability; abla-  
884 tion studies reveal that augmentation choices and random seeds introduce  
885 variance comparable to architectural differences. Furthermore, generalization  
886 experiments on an independent dataset showed substantial ranking inversions  
887 under distribution shift, with foundation models maintaining performance  
888 while classical architectures fail. Moreover, sample size sensitivity analysis  
889 showed that stable rankings require at least 50-100 samples; and within-  
890 distribution training on varying sample sizes reveals that model ranking  
891 hierarchies are dataset-specific.

892 **Author contributions: CRediT**

893 **PMK:** Conceptualization, Data curation, Formal analysis, Methodology,  
894 Software, Visualization, Validation, Writing – original draft; **AAP:** Con-  
895 ceptualization, Data curation, Methodology, Validation, Writing – review  
896 and editing; **YS:** Conceptualization, Methodology, Validation, Writing – re-  
897 view and editing; **ZL:** Conceptualization, Methodology, Validation, Writing  
898 – review and editing; **MT:** Data curation, Funding acquisition, Resources,  
899 Validation, Writing – review and editing; **AC:** Investigation, Data curation,  
900 Resources, Validation, Writing – review and editing; **EAL:** Conceptualization,  
901 Resources, Funding acquisition, Providing samples, Project administration,  
902 Supervision, Writing – review and editing; **SA:** Conceptualization, Funding  
903 acquisition, Methodology, Validation, Supervision, Writing – review and edit-  
904 ing.

905

906 **Declaration of competing interest**

907 The authors declare that they have no known competing financial interests  
908 or personal relationships that could have appeared to influence the work  
909 reported in this paper.

910 **Acknowledgements**

911 The authors gratefully acknowledge the contributions of Brijesh Kumar  
912 Singh, and Priyanka Gupta (Duke-NUS Medical School, Singapore), Roshni  
913 R. Singaraja (National University of Singapore), and Elisa Octavia Velicu  
914 (University of Medicine and Pharmacy Carol Davila, Bucharest) for their  
915 scientific and technical support in this study.

916 **Funding**

917 This work is supported by the Innovation Fund Denmark for the project  
918 DIREC (9142-00001B), CNCS—UEFISCDI, project number PN-III-P4-PCE-  
919 2021-1680. The content is solely the responsibility of the authors and does  
920 not represent the official views of funding sources.

921 **Data availability**

922 Data will be made available on request.

923 **References**

- 924 [1] S. D. Birare, D. V. Swami, S. S. Gaikwad, R. R. Malpani, Study of the  
925 histopathological changes in the heart in postmortem/autopsy cases in  
926 tertiary care center, *Heart Views* 25 (1) (2024) 9–12.
- 927 [2] R. Dinc, Post-stenting restenosis in coronary artery diseases-  
928 pathophysiology and morphological features of re-stenotic tissue, *Biomed-  
929 ical Journal of Scientific & Technical Research* 53 (3) (2023) 44836–44841.
- 930 [3] Y. Matsuhira, R. Shutta, H. Nakamura, K. Yasumoto, K. Yasumura,  
931 A. Tanaka, Y. Matsunaga-Lee, D. Nakamura, M. Yano, M. Yamato, et al.,  
932 Histological findings of rapid progression of neoatherosclerosis including  
933 calcification in hemodialysis patients, *Coronary Artery Disease* 31 (5)  
934 (2020) 479–480.

- 935 [4] R. C. Hubrecht, E. Carter, The 3rs and humane experimental technique:  
936 implementing change, *Animals* 9 (10) (2019) 754.
- 937 [5] R. Arasu, A. Arasu, J. Muller, Carotid artery stenosis: An approach  
938 to its diagnosis and management, *Australian journal of general practice*  
939 50 (11) (2021) 821–825.
- 940 [6] W. C. Roberts, B. P. Everett, V. S. Won, N. Kondapalli, Diagnostic  
941 usefulness of histological examination of the left ventricular “core” excised  
942 to insert a left ventricular assist device in patients with severe heart  
943 failure, *The American Journal of Cardiology* 137 (2020) 71–76.
- 944 [7] G.-R. Pandelea-Dobrovicescu, M. Prodana, F. Golgovici, D. Ionita,  
945 M. Sajin, I. Demetrescu, Surface morphology and histopathological  
946 aspects of metallic used cardiovascular cochr stents, *Metals* 10 (9) (2020)  
947 1112.
- 948 [8] H. M. Ahmad, M. J. Khan, A. Yousaf, S. Ghuffar, K. Khurshid, Deep  
949 learning: a breakthrough in medical imaging, *Current Medical Imaging*  
950 16 (8) (2020) 946–956.
- 951 [9] J. De Matos, S. T. M. Ataky, A. de Souza Britto Jr, L. E. Soares de  
952 Oliveira, A. Lameiras Koerich, Machine learning methods for histopatho-  
953 logical image analysis: A review, *Electronics* 10 (5) (2021) 562.
- 954 [10] B. Kong, Z. Li, S. Zhang, Toward large-scale histopathological image anal-  
955 ysis via deep learning, in: *Biomedical Information Technology*, Elsevier,  
956 2020, pp. 397–414.
- 957 [11] G. Sharnai, A. Livne, A. Polónia, E. Sabo, A. Cretu, G. Bar-Sela, R. Kim-  
958 mel, Deep learning-based image analysis predicts pd-l1 status from h&e-  
959 stained histopathology images in breast cancer, *Nature Communications*  
960 13 (1) (2022) 6753.
- 961 [12] R. Basla, L. Giulivi, L. Magri, G. Boracchi, An expert-driven data  
962 generation pipeline for histological images, in: *2024 IEEE International*  
963 *Symposium on Biomedical Imaging (ISBI)*, IEEE, 2024, pp. 1–5.
- 964 [13] O. Ronneberger, P. Fischer, T. Brox, U-net: Convolutional networks for  
965 biomedical image segmentation, in: *International Conference on Medical*

- 966 image computing and computer-assisted intervention, Springer, 2015, pp.  
967 234–241.
- 968 [14] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, H. Adam, Encoder-  
969 decoder with atrous separable convolution for semantic image segmenta-  
970 tion, in: Proceedings of the European Conference on Computer Vision  
971 (ECCV), 2018, pp. 801–818.
- 972 [15] E. Xie, W. Wang, Z. Yu, A. Anandkumar, J. M. Alvarez, P. Luo,  
973 Segformer: Simple and efficient design for semantic segmentation with  
974 transformers, Advances in neural information processing systems 34  
975 (2021) 12077–12090.
- 976 [16] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson,  
977 T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo, et al., Segment anything,  
978 in: Proceedings of the IEEE/CVF international conference on computer  
979 vision, 2023, pp. 4015–4026.
- 980 [17] J. Ma, Y. He, F. Li, L. Han, C. You, B. Wang, Segment anything in  
981 medical images, Nature Communications 15 (1) (2024) 654.
- 982 [18] S. Yang, J. Feng, X. Mi, H. Bi, H. Zhang, J. Sun, Improved baselines  
983 with synchronized encoding for universal medical image segmentation, in:  
984 Medical Image Computing and Computer-Assisted Intervention (MIC-  
985 CAI), 2025.
- 986 [19] A. F. Cooper, Y. Lu, J. Forde, C. M. De Sa, Hyperparameter optimization  
987 is deceiving us, and how to stop it, Advances in Neural Information  
988 Processing Systems 34 (2021) 3081–3095.
- 989 [20] E. J. Langer, The illusion of control., Journal of personality and social  
990 psychology 32 (2) (1975) 311–328.
- 991 [21] Z. Zhou, M. M. R. Siddiquee, N. Tajbakhsh, J. Liang, Unet++: Redesign-  
992 ing skip connections to exploit multiscale features in image segmentation,  
993 IEEE transactions on medical imaging 39 (6) (2020) 1856–1867.
- 994 [22] M. Z. Alom, C. Yakopcic, M. Hasan, T. M. Taha, V. K. Asari, Recurrent  
995 residual U-Net for medical image segmentation, Journal of Medical  
996 Imaging 6 (1) (2019) 014006.

- 997 [23] S. Paheding, A. A. Reyes-Angulo, M. S. Alam, U-net-based medical image  
998 segmentation: A comparative analysis and future trends, in: *Computer*  
999 *Vision: Challenges, Trends, and Opportunities*, Chapman and Hall/CRC,  
1000 2024, pp. 128–147.
- 1001 [24] H. Yu, L. T. Yang, Q. Zhang, D. Armstrong, M. J. Deen, Convolutional  
1002 neural networks for medical image analysis: state-of-the-art, comparisons,  
1003 improvement and perspectives, *Neurocomputing* 444 (2021) 92–110.
- 1004 [25] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez,  
1005 Ł. Kaiser, I. Polosukhin, Attention is all you need, *Advances in neural*  
1006 *information processing systems* 30 (2017).
- 1007 [26] G. Chrysos, Architecture design: From neural networks to foundation  
1008 models, in: *2024 IEEE 11th International Conference on Data Science*  
1009 *and Advanced Analytics (DSAA)*, IEEE, 2024, pp. 1–3.
- 1010 [27] K. Xia, J. Wang, Recent advances of transformers in medical image  
1011 analysis: a comprehensive review, *MedComm–Future Medicine* 2 (1)  
1012 (2023) e38.
- 1013 [28] Y. Lai, Application and effectiveness evaluation of bayesian optimization  
1014 algorithm in hyperparameter tuning of machine learning models, in: *2024*  
1015 *International Conference on Power, Electrical Engineering, Electronics*  
1016 *and Control (PEEEEC)*, IEEE, 2024, pp. 351–355.
- 1017 [29] N. Subaşı, Comprehensive analysis of grid and randomized search on  
1018 dataset performance, *European Journal of Engineering and Applied*  
1019 *Sciences* 7 (2) (2024) 77–83.
- 1020 [30] Monica, P. Agrawal, A survey on hyperparameter optimization of machine  
1021 learning models, in: *2024 2nd International Conference on Disruptive*  
1022 *Technologies (ICDT)*, IEEE, 2024, pp. 11–15.
- 1023 [31] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang,  
1024 W. Chen, LoRA: Low-rank adaptation of large language models, in:  
1025 *International Conference on Learning Representations*, 2022.
- 1026 [32] A. W. Salehi, S. Khan, G. Gupta, B. I. Alabdullah, A. Almjally, H. Al-  
1027 solai, T. Siddiqui, A. Mellit, A study of cnn and transfer learning in

- 1028 medical imaging: Advantages, challenges, future scope, Sustainability  
1029 15 (7) (2023) 5930.
- 1030 [33] G.-S. Hong, M. Jang, S. Kyung, K. Cho, J. Jeong, G. Y. Lee, K. Shin,  
1031 K. D. Kim, S. M. Ryu, J. B. Seo, et al., Overcoming the challenges in the  
1032 development and implementation of artificial intelligence in radiology: a  
1033 comprehensive review of solutions beyond supervised learning, Korean  
1034 Journal of Radiology 24 (11) (2023) 1061–1080.
- 1035 [34] L. Maier-Hein, M. Eisenmann, A. Reinke, S. Onogur, M. Stankovic,  
1036 P. Scholz, T. Arbel, H. Bogunovic, A. P. Bradley, A. Carass, et al., Why  
1037 rankings of biomedical image analysis competitions should be interpreted  
1038 with care, Nature communications 9 (1) (2018) 5217.
- 1039 [35] J. Lin, The neural hype and comparisons against weak baselines, in:  
1040 Acm sigir forum, Vol. 52, ACM New York, NY, USA, 2018, pp. 40–51.
- 1041 [36] J. Lin, D. Campos, N. Craswell, B. Mitra, E. Yilmaz, Significant improve-  
1042 ments over the state of the art? a case study of the ms marco document  
1043 ranking leaderboard, in: Proceedings of the 44th international ACM  
1044 SIGIR conference on research and development in information retrieval,  
1045 2021, pp. 2283–2287.
- 1046 [37] H. Ranglani, Empirical analysis of the bias-variance tradeoff across  
1047 machine learning models, Machine Learning and Applications: An Inter-  
1048 national Journal (MLAIJ) 11 (4) (2024).
- 1049 [38] A. Makarova, H. Shen, V. Perrone, A. Klein, J. B. Faddoul, A. Krause,  
1050 M. Seeger, C. Archambeau, Overfitting in bayesian optimization: an  
1051 empirical study and early-stopping solution, in: 2nd Workshop on Neural  
1052 Architecture Search (NAS 2021)@ ICLR 2021, NAS 2021, 2021.
- 1053 [39] L. Hertel, P. Baldi, D. L. Gillen, Reproducible hyperparameter optimiza-  
1054 tion, Journal of Computational and Graphical Statistics 31 (1) (2022)  
1055 84–99.
- 1056 [40] I. V. Tetko, R. van Deursen, G. Godin, Be aware of overfitting by  
1057 hyperparameter optimization!, Journal of Cheminformatics 16 (1) (2024)  
1058 139.

- 1059 [41] P. Xu, X. Ji, M. Li, W. Lu, Small data machine learning in materials  
1060 science, *npj Computational Materials* 9 (1) (2023) 42.
- 1061 [42] D. Patil, Explainable artificial intelligence (xai): Enhancing transparency  
1062 and trust in machine learning models, Available at SSRN 5057400 (2024).
- 1063 [43] M. A. Islam, M. A. Jahin, M. Mridha, N. Dey, A unified framework for  
1064 evaluating the effectiveness and enhancing the transparency of explainable  
1065 ai methods in real-world applications, arXiv preprint arXiv:2412.03884  
1066 (2024).
- 1067 [44] A. Curaj, Z. Wu, M. Staudt, E. A. Liehn, Induction of accelerated  
1068 atherosclerosis in mice: The “wire-injury” model, *Journal of Visualized  
1069 Experiments* (162) (2020). doi:10.3791/54571.
- 1070 [45] Z. Liu, Q. Lv, C. H. Lee, L. Shen, Segmenting medical images with  
1071 limited data, *Neural Networks* 177 (2024) 106367.
- 1072 [46] J. Xu, M. Li, Z. Zhu, Automatic data augmentation for 3d medical image  
1073 segmentation, in: *International Conference on Medical Image Computing  
1074 and Computer-Assisted Intervention*, Springer, 2020, pp. 378–387.
- 1075 [47] V. R. Dasari, B. E. Geerhart, D. M. Alexander, Increasing image seg-  
1076 mentation accuracy on small datasets by merging multiple inferences on  
1077 augmented images, in: *Disruptive Technologies in Information Sciences  
1078 VI*, Vol. 12117, SPIE, 2022, p. 1211709.
- 1079 [48] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image  
1080 recognition, in: *Proceedings of the IEEE conference on computer vision  
1081 and pattern recognition*, 2016, pp. 770–778.
- 1082 [49] M. Tan, Q. Le, Efficientnet: Rethinking model scaling for convolutional  
1083 neural networks, in: *International conference on machine learning*, PMLR,  
1084 2019, pp. 6105–6114.
- 1085 [50] G. Huang, Z. Liu, L. Van Der Maaten, K. Q. Weinberger, Densely  
1086 connected convolutional networks, in: *Proceedings of the IEEE conference  
1087 on computer vision and pattern recognition*, 2017, pp. 4700–4708.
- 1088 [51] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand,  
1089 M. Andreetto, H. Adam, Mobilenets: Efficient convolutional neural

- 1090 networks for mobile vision applications, arXiv preprint arXiv:1704.04861  
1091 (2017).
- 1092 [52] E. Khalili, B. Priego-Torres, A. Leon-Jimenez, D. Sanchez-Morillo, Auto-  
1093 matic lung segmentation in chest x-ray images using sam with prompts  
1094 from yolo, *IEEE Access* 12 (2024) 122805–122819.
- 1095 [53] K. Xu, L. Goetz, N. Rajpoot, On generalisability of segment anything  
1096 model for nuclear instance segmentation in histology images, arXiv  
1097 preprint arXiv:2401.14248 (2024).
- 1098 [54] J. Bergstra, R. Bardenet, Y. Bengio, B. Kégl, Algorithms for hyper-  
1099 parameter optimization, in: *Advances in Neural Information Processing*  
1100 *Systems*, Vol. 24, 2011.
- 1101 [55] H. Wang, K. Yang, Bayesian optimization, in: *Many-Criteria Optimiza-*  
1102 *tion and Decision Analysis: State-of-the-Art, Present Challenges, and*  
1103 *Future Perspectives*, Springer, 2023, pp. 271–297.
- 1104 [56] M. Yeung, L. Rundo, Y. Nan, E. Sala, C.-B. Schönlieb, G. Yang, Calibrat-  
1105 ing the dice loss to handle neural network overconfidence for biomedical  
1106 image segmentation, *Journal of Digital Imaging* 36 (2) (2023) 739–752.
- 1107 [57] D. Müller, I. Soto-Rey, F. Kramer, Towards a guideline for evaluation  
1108 metrics in medical image segmentation, *BMC Research Notes* 15 (1)  
1109 (2022) 210.
- 1110 [58] M. Bhagat, B. Bakariya, A comprehensive review of cross-validation  
1111 techniques in machine learning, *IJSAT-International Journal on Science*  
1112 *and Technology* 16 (1) (2025).
- 1113 [59] M. Yeung, E. Sala, C.-B. Schönlieb, L. Rundo, Unified focal loss: Gener-  
1114 alising dice and cross entropy-based losses to handle class imbalanced  
1115 medical image segmentation, *Computerized Medical Imaging and Graph-*  
1116 *ics* 95 (2022) 102026.
- 1117 [60] B. Efron, Bootstrap methods: another look at the jackknife, in: *Break-*  
1118 *throughs in statistics: Methodology and distribution*, Springer, 1992, pp.  
1119 569–593.

- 1120 [61] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion,  
1121 O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, et al.,  
1122 Scikit-learn: Machine learning in python, the Journal of machine Learning  
1123 research 12 (2011) 2825–2830.
- 1124 [62] M. Friedman, The use of ranks to avoid the assumption of normality  
1125 implicit in the analysis of variance, Journal of the american statistical  
1126 association 32 (200) (1937) 675–701.
- 1127 [63] J. Demšar, Statistical comparisons of classifiers over multiple data sets,  
1128 Journal of Machine learning research 7 (Jan) (2006) 1–30.
- 1129 [64] P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy,  
1130 D. Cournapeau, E. Burovski, P. Peterson, W. Weckesser, J. Bright, et al.,  
1131 Scipy 1.0: fundamental algorithms for scientific computing in python,  
1132 Nature methods 17 (3) (2020) 261–272.
- 1133 [65] J. Cohen, Statistical power analysis for the behavioral sciences, 2nd  
1134 Edition, Routledge, 1988.
- 1135 [66] D. Karimi, S. E. Salcudean, Reducing the hausdorff distance in medical  
1136 image segmentation with convolutional neural networks, IEEE Transactions  
1137 on Medical Imaging 39 (2) (2020) 499–513.
- 1138 [67] A. Reinke, M. D. Tizabi, M. Baumgartner, et al., Understanding metric-  
1139 related pitfalls in image analysis validation, Nature Methods 21 (2024)  
1140 182–194.
- 1141 [68] Hugging Face, Fine-tune a semantic segmentation model with a cus-  
1142 tom dataset, <https://huggingface.co/blog/fine-tune-segformer>,  
1143 accessed: 2026-02-07 (2022).
- 1144 [69] H. Gu, H. Dong, J. Yang, M. A. Mazurowski, How to build the best med-  
1145 ical image segmentation algorithm using foundation models: a compre-  
1146 hensive empirical study with segment anything model, Machine Learning  
1147 for Biomedical Imaging (MELBA) 3 (2025) 88–120.
- 1148 [70] I. Loshchilov, F. Hutter, Decoupled weight decay regularization, in:  
1149 International Conference on Learning Representations, 2019.

1150 **Supplementary Materials**

1151 **Appendix A. Bayesian Hyperparameter Optimization Details**

1152 Our hyperparameter optimization uses the Tree-structured Parzen Estimator (TPE) [54], implemented via Optuna with a fixed random seed ( $seed = 42$ )  
1153 and 10 startup trials of random sampling before switching to the TPE sampler.  
1154

1155 **Why TPE, not Gaussian Processes.** Classical Bayesian optimization  
1156 fits a Gaussian Process (GP) to model the objective  $f(x)$  directly, then  
1157 maximizes an acquisition function (e.g., Expected Improvement) over that  
1158 surrogate. This requires inverting a covariance matrix at cost  $\mathcal{O}(t^3)$  per step,  
1159 which becomes prohibitive as the number of trials  $t$  grows. TPE avoids this  
1160 entirely by modeling the *inverse* relationship: instead of  $p(y | x)$ , it models  
1161  $p(x | y)$ .

1162 **Observation splitting.** After  $t$  evaluations, TPE partitions the observed  
1163 hyperparameter–score pairs at a quantile threshold  $\gamma$ . Configurations with  
1164 scores better than the threshold  $y^*$  define a “good” density  $\ell(x)$ ; the remainder  
1165 define a “bad” density  $g(x)$ :

$$p(x | y) = \begin{cases} \ell(x) & \text{if } y < y^* \\ g(x) & \text{if } y \geq y^* \end{cases} \quad (\text{S1})$$

1166 where  $y^*$  is the  $\gamma$ -quantile of observed scores ( $\gamma = 0.25$  in Optuna) and both  
1167  $\ell(x)$  and  $g(x)$  are estimated via Parzen window (kernel density) estimators.

1168 **Selection criterion.** Bergstra et al. [54] show that, under this formulation,  
1169 the next configuration should be chosen to maximize the ratio of the two  
1170 densities:

$$x_{\text{next}} = \arg \max_x \frac{\ell(x)}{g(x)} \quad (\text{S2})$$

1171 Intuitively, this selects hyperparameters that are *likely under good trials* and  
1172 *unlikely under bad trials*. In practice, TPE draws many candidates from  
1173  $\ell(x)$  and picks the one with the highest  $\ell(x)/g(x)$  ratio. This naturally  
1174 balances exploitation (sampling where good configurations cluster) with  
1175 exploration (avoiding regions dominated by poor configurations), enabling  
1176 efficient navigation of high-dimensional hyperparameter spaces.

1177 **Appendix B. Hardware and Software Specifications**

1178 For complete reproducibility, we provide detailed specifications of our  
1179 computational environment in Table S14.

Table S14: Hardware and Software Environment Specifications

<b>Component</b>	<b>Specification</b>
<i><b>Hardware (Main Experiments)</b></i>	
GPU	1–3× NVIDIA Tesla V100 (32GB VRAM each)
GPU Architecture	Volta
CUDA Cores	5,120 per GPU
Memory Bandwidth	900 GB/s
<i><b>Hardware (Resolution Ablation)</b></i>	
GPU	3× NVIDIA Tesla V100-SXM2 (32GB VRAM each)
GPU Architecture	Volta
Provider	UCloud ( <a href="https://cloud.sdu.dk">https://cloud.sdu.dk</a> )
<i><b>Software Environment</b></i>	
Deep Learning Framework	PyTorch 2.1+
CUDA Version	12.1
cuDNN Version	8.9
Python Version	3.10+
Segmentation Library	segmentation-models-pytorch 0.3.3+
Foundation Models	transformers 4.36+
Hyperparameter Optimization	Optuna 3.4+
Experiment Tracking	JSON-based offline storage
<i><b>Reproducibility Settings</b></i>	
Random Seed	42
Python Random	Seeded
NumPy Random	Seeded
PyTorch Manual Seed	Seeded
CUDA Manual Seed	Seeded (all devices)
cuDNN Deterministic	True
cuDNN Benchmark	False
Mixed Precision (AMP)	Enabled
<i><b>Training Configuration</b></i>	
Data Workers	8
Pin Memory	Enabled
Early Stopping Patience	20 epochs
Maximum Epochs	200

1180 Each individual experiment was executed on a single NVIDIA Tesla V100  
1181 GPU to eliminate hardware-induced variance; up to three identical GPUs  
1182 were used concurrently to run different model experiments in parallel. The  
1183 deterministic settings ensure that given the same random seed, data splits,  
1184 and hyperparameters, our results can be practically reproduced. We note that  
1185 Automatic Mixed Precision (AMP) may introduce minor non-determinism in  
1186 floating-point reduction operations; however, this has negligible impact on  
1187 aggregate metrics and conclusions.