

Noise-Aware Dynamic Convolution for Improved Generalizability of Retinal Disease Diagnosis using Optical Coherence Tomography Images

Deeksha Chutani^a, Phaneendra K. Yalavarthy^{a,b,*}

^aDepartment of Computational and Data Sciences, Indian Institute of Science, Bangalore, Karnataka, 560012, India

^bTANUH-AI Centre of Excellence in Healthcare, Indian Institute of Science, Bangalore, Karnataka, 560012, India

Abstract.

Significance: Optical coherence tomography (OCT) is widely used for the diagnosis of retinal diseases. However, deep learning models trained on a single dataset often degrade when deployed across scanners and clinical sites due to device-dependent speckle variability and acquisition differences, limiting their reliability in real-world screening.

Aim: This study aims to develop a lightweight deep learning framework that leverages speckle characteristics in optical coherence tomography (OCT) images to improve cross-scanner generalizability for retinal disease classification while preserving real-time inference efficiency.

Approach: We propose NA-DyCNN, a noise-aware dynamic convolutional neural network that minimizes the expected classification risk over multiple stochastic realizations of multiplicative speckle perturbations and regularizes the dynamic routing mechanism to produce scanner-invariant kernel mixtures. The framework was evaluated using over 105,000 B-scans from three heterogeneous OCT cohorts under strict zero-shot cross-dataset transfer.

Results: NA-DyCNN consistently outperformed lightweight baselines across four zero-shot cross-dataset transfer scenarios, achieving up to 92.87% accuracy, a weighted F_2 score of 92.89%, and Cohen's κ of up to 0.896, demonstrating improved robustness and generalization under cross-dataset shifts. The model maintained high efficiency with only 0.4M parameters and an inference latency of 0.53 ms per B-scan on an NVIDIA GB10 GPU.

Conclusions: Modeling post-acquisition speckle variability during training improves the generalizability of OCT classifiers without increasing the inference cost, thereby enabling the more reliable deployment of artificial intelligence (AI)-assisted retinal screening across heterogeneous imaging systems.

Keywords: Optical Coherence Tomography, Deep Learning, Generalizability, Dynamic Convolution, Retinal Disease Classification, Conformal Prediction.

*Phaneendra K. Yalavarthy, yalavarthy@iisc.ac.in

1 Introduction

Optical coherence tomography (OCT) is an optical imaging modality based on low-coherence interferometry that enables noninvasive micrometer-resolution visualization of retinal microstructures, supporting the diagnosis and monitoring of diseases, such as age-related macular degeneration and diabetic macular edema.^{1,2} A fundamental characteristic of OCT imaging is the presence of speckle patterns arising from the coherent interference of backscattered light from sub-resolution tissue scatterers. These granular intensity fluctuations occur due to constructive and

35 destructive interference of scattered optical waves and are therefore intrinsic to all coherent imag-
36 ing systems. Importantly, OCT speckle exhibits a dual role: while it reduces visual image quality
37 and contrast, its statistical properties encode information about tissue microstructure and optical
38 scattering characteristics.³ Consequently, speckle patterns are increasingly recognized not merely
39 as noise but as signal-carrying features that may contain diagnostically relevant information.

40 The increasing global burden of retinal disease has motivated the widespread adoption of deep
41 learning for automated OCT analysis, with convolutional neural networks (CNNs),⁴ transfer learn-
42 ing approaches,⁵ and Vision Transformer-based models⁶ achieving strong performance on curated
43 OCT benchmarks.⁷

44 Despite these advances, most reported results have been obtained under single-dataset condi-
45 tions and fail to reliably translate to new scanners and clinical sites. In real-world deployments,
46 OCT systems differ in terms of optical design, coherence length, detector sensitivity, beam scan-
47 ning strategies, and reconstruction pipelines. These factors influence the statistical characteris-
48 tics of speckle patterns and image contrast, leading to systematic variations in OCT appearance
49 across devices. Consequently, models trained on a single acquisition regime frequently overfit
50 scanner-specific texture cues and suffer substantial performance degradation under cross-site or
51 cross-scanner transfer.^{8,9}

52 A defining yet underutilized contributor to this domain shift is the scanner-dependent variability
53 of the speckle patterns present in clinically archived OCT images. Speckle formation depends
54 on factors, such as illumination coherence, optical beam properties, detector aperture, and tissue
55 microstructure.³ Because these parameters vary across imaging systems, speckle statistics differ
56 between scanners, even for anatomically similar retinas.

57 Historically, most OCT image processing research has treated speckle primarily as a source of

58 signal degradation, motivating extensive work on denoising and speckle suppression through filter-
59 ing, averaging, or wavelet-based methods.¹⁰⁻¹² However, suppressing speckle may also remove the
60 acquisition-dependent information that characterizes the imaging system. Paluru et al. proposed
61 a self-distillation framework to improve robustness;¹³ nevertheless, such methods remain agnostic
62 to the physical origins of speckle formation and do not explicitly model scanner-dependent noise
63 processes.

64 In contrast, the practical deployment of OCT classifiers requires models that are both compu-
65 tationally efficient and clinically reliable. Lightweight architectures, such as MobileNet variants,¹⁴
66 reduce inference cost but often exhibit fragile generalization under dataset shift, whereas high-
67 capacity networks, such as ResNet¹⁵ are incompatible with point-of-care constraints.¹⁶ Moreover,
68 clinical translation demands more than raw accuracy; calibrated uncertainty estimates and safe
69 failure mechanisms are essential for trustworthy screening.

70 Dynamic convolution provides a promising compromise between expressivity and efficiency by
71 adapting convolutional kernels to individual inputs.¹⁷ However, its potential to explicitly account
72 for scanner-induced speckle variability in OCT imaging has not been investigated in prior work.

73 In this study, we introduce a Noise-Aware Dynamic Convolutional Neural Network (NA-
74 DyCNN) for robust retinal OCT classification under scanner-induced shifts. The key contributions
75 of this work are summarized as follows:

- 76 • **Noise-aware dynamic convolution framework.** We propose NA-DyCNN, a dynamic con-
77 volutional architecture that integrates a speckle-variability-aware training objective. By opti-
78 mizing the expected classification risk over multiple stochastic realizations of post-acquisition
79 speckle perturbations, the model learns routing functions that are robust to scanner-dependent

80 noise characteristics while preserving inference-time efficiency.

81 • **Cross-dataset robustness through noise-aware training.** Unlike conventional data aug-
82 mentation approaches, the proposed objective explicitly regularizes the dynamic routing
83 mechanism by enforcing consistent kernel mixtures across multiple noise realizations of
84 the same OCT B-scan, reducing reliance on scanner-specific texture cues.

85 • **Large-scale cross-dataset evaluation under strict zero-shot transfer.** We conduct a com-
86 prehensive evaluation across more than 105,000 OCT B-scans from three heterogeneous
87 clinical cohorts (Table 1), demonstrating improved generalization across cross-scanner and
88 cross-site domain shifts compared with static CNNs, dynamic CNNs, and lightweight mobile
89 architectures.

90 • **Clinical reliability and safety analysis.** Beyond diagnostic accuracy, we evaluate cali-
91 bration quality, Cohen’s κ , recall-weighted F_2 score, and uncertainty attribution using the
92 Error–Uncertainty (Err.Unc) metric. We further assessed deployment safety using Mondrian
93 conformal prediction, providing finite-sample coverage guarantees under dataset shift.

94 **2 Methods and Materials**

95 *2.1 Datasets*

Table 1 OCT dataset characteristics and usage protocols after quality control (QC).

| Dataset | Source / Scanner | Images | Classes Used | Resolution Profile (post-QC) | Usage Protocol |
|----------------|--|--------|---------------------------------------|--|--|
| D_1 – UCSD | Heidelberg Spectralis SD-OCT (Multi-center) ⁷ | ~84k | CNV, DME, DRUSEN, NOR-MAL | 512×496 (60%), 768×496 (31%), 512×512 (9%) | In-dataset training; cross-dataset test for $D_2 \rightarrow D_1$ |
| D_2 – OCT-C8 | Curated multi-vendor OCT ¹⁸ | ~10.9k | CNV, DME, DRUSEN, NORMAL [†] | 512×496 (~48%), 512×512 (~27%), 768×496 (~25%) | In-dataset training; cross-dataset test for $D_1 \rightarrow D_2$ |
| D_3 – NEH | Heidelberg Spectralis SD-OCT (Noor Eye Hospital, Tehran, Iran) ¹⁹ | ~10.9k | CNV, DRUSEN, NORMAL | Uniform 768×496 (100%) | External test only ($D_1 \rightarrow D_3$, $D_2 \rightarrow D_3$) |

[†] Four classes selected from original eight to ensure cross-dataset consistency.

96 We evaluated the proposed framework across three large, distinct retinal OCT datasets acquired
 97 using heterogeneous scanner hardware and clinical acquisition protocols. To enable robust cross-
 98 dataset evaluation, we restrict our analysis to four clinically significant classes: Choroidal Neovas-
 99 cularization (CNV), Diabetic Macular Edema (DME), Drusen, and Normal. A detailed summary
 100 of the dataset characteristics is provided in Table 1.

101 D_1 (UCSD)⁷ was split at the patient level²⁰ into 70/15/15 train/validation/test subsets after
 102 filtering rare or highly anisotropic resolutions, retaining only images of size 512 × 496, 768 × 496,
 103 and 512 × 512. This dataset was curated from multiple clinical centers and labeled through a multi-
 104 tier grading pipeline involving trained graders, ophthalmologists, and senior retinal specialists.⁷

105 D_2 (OCT-C8)¹⁸ served as an independent multi-vendor cohort, where only four overlapping
 106 classes were retained and the original dataset splits were preserved after identical resolution fil-
 107 tering. This dataset aggregates retinal OCT images from multiple sources and therefore reflects

108 heterogeneous acquisition environments and scanner configurations.¹⁸

109 In D_3 (NEH),¹⁹ all images were resolution-uniform (768×496), and only B-scans whose
110 frame-level labels matched clinician-assigned volume labels were retained to construct a patient-
111 balanced evaluation cohort across CNV, Drusen, and Normal classes. The NEH dataset contains
112 more than 16,000 OCT B-scans acquired at Noor Eye Hospital (Tehran, Iran) and includes frame-
113 level labels assigned by a retinal specialist.¹⁹ The publicly available NEH dataset does not include
114 DME annotations; therefore, all evaluations involving D_3 were performed in a three-class setting
115 using only overlapping diagnostic categories.

116 Two independent training configurations were evaluated to assess cross-dataset generalization.
117 In the first configuration, models were trained on D_1 (UCSD) and evaluated on both D_2 (OCT-C8)
118 and D_3 (NEH). In the second configuration, models were trained on D_2 (OCT-C8) and evaluated
119 on D_1 (UCSD) and D_3 (NEH). Since D_3 (NEH) does not include DME annotations, models trained
120 on D_3 would produce a three-class classifier incompatible with the four-class evaluation protocol
121 of D_1 and D_2 . Consequently, D_3 serves exclusively as a target domain in all transfer evaluations.

122 Although D_1 (UCSD) and D_3 (NEH) share the Heidelberg Spectralis SD-OCT platform, they
123 differ substantially in acquisition protocols, clinical sites, patient demographics, and annotation
124 practices. Consequently, $D_1 \rightarrow D_3$ represents a cross-site shift without a scanner change, whereas
125 all D_2 transfers involve both cross-scanner and cross-site variability.

126 All datasets consist of vendor-exported clinical B-scans stored in standard image formats. D_1
127 (UCSD) contains JPEG images, D_2 (OCT-C8) contains JPG images, and D_3 (NEH) contains JPG
128 and TIFF images. These formats reflect typical post-processed exports used in hospital PACS
129 systems and telemedicine transmission.

130 Together, these datasets represent heterogeneous clinical cohorts acquired across different in-

stitutions, annotation pipelines, and imaging environments, providing a suitable benchmark for evaluating cross-dataset robustness of OCT classification models.

2.2 Image Preprocessing and Data Augmentation

Raw B-scans were standardized to 8-bit grayscale, and 16-bit inputs were linearly scaled to preserve the luminance structures. To maintain the retinal layer morphology and device-specific texture characteristics, the images were square-padded to $\max(W, H)$ before being downsampled to 256×256 pixels using Lanczos interpolation.²¹ Finally, per-image Z-normalization was applied to ensure numerical stability and contrast consistency across heterogeneous scanner hardware. During training, mild spatial (horizontal flip, rotation, affine translation, and scaling) and photometric (brightness and contrast jitter) augmentations were applied; no augmentation was used during the evaluation. All experiments were implemented in PyTorch²² using deterministic settings and three random seeds as a computational trade-off across multiple cross-dataset transfer settings, ablation studies, calibration analysis, and uncertainty experiments. To verify that three seeds were sufficient, we additionally evaluated the highest-variance transfer setting (Dynamic CNN on $D_1 \rightarrow D_3$) using five seeds. The resulting performance (F_2 score) changed only marginally from 66.78 ± 8.9 (three seeds) to 65.37 ± 8.7 (five seeds), indicating that the observed variability was primarily due to task difficulty rather than insufficient seed sampling. Importantly, the relative ranking of the methods remained unchanged, with NA-DyCNN continuing to outperform Dynamic CNN by a substantial margin on this challenging transfer setting (72.54 ± 2.1 vs. 65.37 ± 8.7), suggesting that the main conclusions of the study are robust to additional seed sampling.

The OCT datasets (Table 1) used in this study contained only vendor-exported, log-compressed B-scans. Raw interferometric or linear-intensity OCT volumes required for first-principles speckle

153 simulations are not accessible in UCSD, OCT-C8, or NEH, which constrains all learning-based
 154 methods to operate in the post-acquisition image domain. Therefore, our objective was not to
 155 reproduce physical speckle formation but to model the residual speckle variability that persists
 156 after vendor-specific post-processing and defines the cross-scanner domain shift in real clinical
 157 exports.

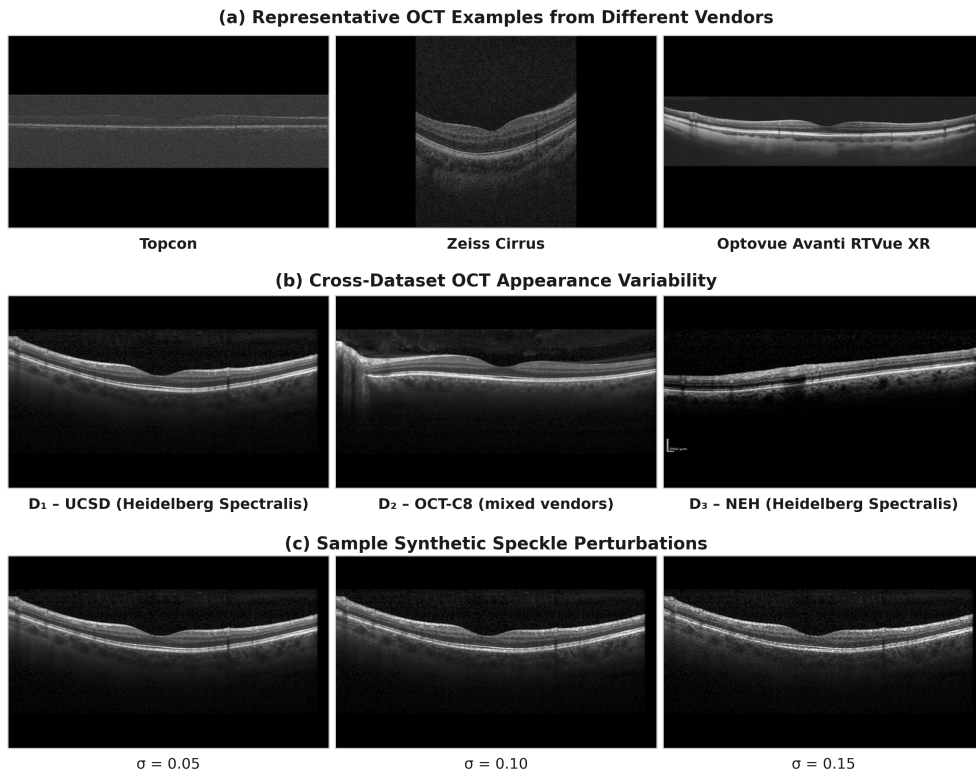


Fig 1 Visualization of OCT appearance variability, together with sample synthetic perturbations used during noise-aware training. (a) Representative OCT B-scans obtained from publicly available sources acquired using Topcon,²³ Zeiss Cirrus,²⁴ and Optovue Avanti RTVue XR²⁵ systems, illustrating the diversity of image appearance, contrast characteristics, and retinal layer representation across different vendors. (b) Representative OCT B-scans from the three datasets used in this study: D_1 (UCSD; Heidelberg Spectralis), D_2 (OCT-C8; mixed vendors), and D_3 (NEH; Heidelberg Spectralis), illustrating scanner-dependent variations in the speckle texture and reflectivity statistics. (c) Noise-perturbed realizations of the same OCT B-scan generated using multiplicative Gaussian perturbations ($\sigma = 0.05, 0.10, \text{ and } 0.15$), illustrating increasing levels of synthetic speckle variability. These values are shown for visualization purpose and do not correspond exactly to the hyperparameters used during training.

158 To visually illustrate the cross-scanner appearance variability motivating our approach, Fig. 1
 159 shows representative scans from different vendors in the first row followed by normal retinal OCT

160 B-scans from the three datasets used in this study in the second row. Noticeable differences in
161 speckle texture, background statistics, and contrast are observed across datasets due to variations
162 in scanner hardware, acquisition parameters, and vendor-specific reconstruction pipelines. For
163 clarity, the third row shows noise-perturbed realizations of the same UCSD normal B-scan us-
164 ing multiplicative perturbations with increasing variance, illustrating the type of residual speckle
165 variability modeled during noise-aware training.

166 *2.3 Baseline Methods*

167 To validate the efficiency and robustness of our framework, we compared it with four distinct
168 architectures. The three custom convolutional designs (Basic CNN, Dynamic CNN, and proposed
169 NA-DyCNN) are illustrated in Fig. 2.

170 All models were trained from scratch to avoid the confounding scanner-domain bias introduced
171 by ImageNet pretraining, which has been shown to encode non-medical texture priors and can
172 artificially inflate cross-domain performance.

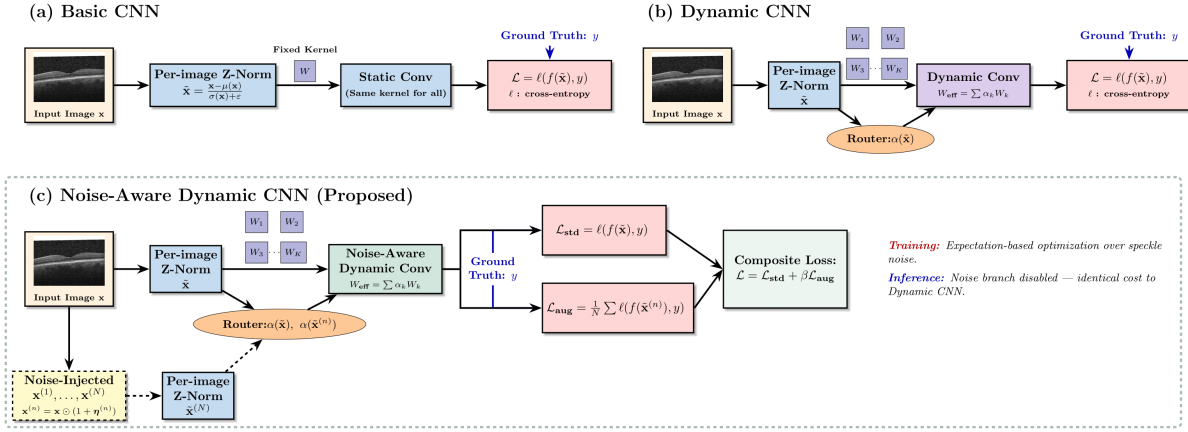


Fig 2 Structural comparison of convolutional architectures for OCT classification. (a) Basic CNN applies per-image z -normalization, followed by static convolution with shared kernels across all inputs. (b) Dynamic CNN employs a router network to generate image-dependent kernel mixtures for dynamic convolution. (c) Noise-Aware Dynamic CNN (Proposed) introduces expectation-based regularization by training on multiple stochastic realizations of scanner-export speckle variability $\mathbf{x}^{(n)}$ of each input, while using the same dynamic routing mechanism for both clean and noisy samples. During inference, the noise branch is disabled, yielding a computational cost identical to that of the standard Dynamic CNN. Here, $\ell(\cdot, \cdot)$ denotes the multi-class cross-entropy loss.

- 173 • **Basic CNN:** A standard hierarchical baseline comprising four sequential convolutional blocks
- 174 (16 \rightarrow 32 \rightarrow 64 \rightarrow 128 filters) followed by Global Average Pooling (GAP) and a two-layer
- 175 MLP classifier (Fig. 2a). This serves as a reference for static feature extraction.
- 176 • **Dynamic CNN:** Based on Chen et al.,¹⁷ this model shares the identical macro-architecture
- 177 and parameter count as the Basic CNN but replaces static layers with dynamic convolution
- 178 modules (with $K = 4$ experts; Fig. 2b). This isolates the specific contribution of the input-
- 179 adaptive processing. Following prior implementations of Dynamic Convolution,¹⁷ which
- 180 commonly use multiple parallel kernels per layer, we set the number of experts to four for
- 181 all dynamic convolution models. This configuration provides a practical trade-off between
- 182 the representational capacity and computational cost.
- 183 • **MobileNetV3-Small:** A state-of-the-art lightweight architecture¹⁴ utilizing depthwise-separable

184 convolutions, included to benchmark performance against established mobile-optimized net-
 185 works.

- 186 • **ResNet18 (reference only)**: A high-capacity residual network¹⁵ was included solely as a
 187 computationally unconstrained performance reference. It was reported only to contextualize
 188 the robustness–efficiency trade-off achieved by lightweight models.

Table 2 Computational complexity and inference efficiency for 256×256 grayscale input on NVIDIA GB 10 with batch size 1. Latency measurements represent the mean \pm std over 1000 forward passes. The relative training cost is normalized to Dynamic CNN (1.0 \times), measured on the D_2 training split (batch size 64). Bold entries denote the proposed method.

| Model | Params (M) | FLOPs (M) | Latency (ms) | FPS | Relative Training Cost | Convolution Type |
|--|-------------|--------------|-----------------------------------|-------------|--------------------------------|-------------------------|
| Basic CNN | 0.11 | 243.8 | 0.25 ± 0.04 | 4000 | 0.90 \times | Static |
| Dynamic CNN | 0.40 | 482.5 | 0.53 ± 0.05 | 1886 | 1.00 \times | Dynamic |
| Proposed (NA-DyCNN)[†] | 0.40 | 482.5 | 0.53 ± 0.05 | 1886 | 4.69\times | Dynamic |
| MobileNetV3 | 1.52 | 75.2 | 1.41 ± 0.06 | 709 | 0.79 \times | Static (Depthwise-sep.) |
| ResNet18 | 11.17 | 2279.0 | 1.26 ± 0.10 | 793 | 1.73 \times | Static (Residual) |

[†]Noise-aware training affects only the training phase; inference cost is identical to Dynamic CNN.

FLOPs are reported for $B=1$ and include router overhead for Dynamic CNN variants.

Relative Training Cost is normalized to Dynamic CNN (5.57 \pm 0.17 s/epoch); For NA-DyCNN, the training time is 26.13 \pm 0.29 s/epoch ($N=4$ realizations).

189 Table 2 summarizes the parameter counts, FLOPs, empirical inference latency, and relative
 190 training cost. As illustrated in Fig. 2c, NA-DyCNN introduces speckle-variability-aware regular-
 191 ization only during training; consequently, its inference-time computational cost is identical to that
 192 of the Dynamic CNN. The training overhead of 4.69 \times relative to Dynamic CNN reflects the cost
 193 of processing $N = 4$ additional noisy forward passes per mini-batch and is incurred only during
 194 training. Because model training is performed offline once and the resulting model is deployed
 195 without the noise branch, this additional cost does not affect the real-world inference efficiency.

196 2.4 Proposed Approach: Noise-Aware Dynamic Convolutional Neural Network

197 We propose a Noise-Aware Dynamic Convolutional Neural Network (NA-DyCNN) (Fig. 2) that
 198 explicitly models the post-acquisition speckle variability present in clinically archived OCT im-

199 ages within the learning objective. Rather than optimizing on a single clean realization, the model
200 minimizes the expected classification risk over a distribution of speckle-corrupted observations,
201 thereby learning features and routing functions that are invariant to scanner-dependent noise statis-
202 tics while preserving real-time inference efficiency.

203 *2.4.1 Motivation: Post-Acquisition Speckle Variability as a Source of Domain Shift*

204 Retinal OCT is a coherent imaging modality^{1,2} and is inherently affected by multiplicative speckle
205 noise arising from the interference of backscattered light from sub-resolution tissue structures.
206 The statistical properties of this noise vary across scanner hardware due to differences in coher-
207 ence length, detector sensitivity, and signal averaging strategies, making vendor-processed speckle
208 variability a dominant observable source of cross-dataset shift in clinical exports. Conventional
209 and Dynamic CNNs trained on single image realizations tend to overfit scanner-specific noise sig-
210 natures, degrading generalization to unseen acquisition systems.

211 *2.4.2 Realistic Speckle Modeling in Processed Clinical Data*

212 While theoretical OCT speckle statistics are best modeled using raw 12-bit or 16-bit intensity data,
213 standard machine learning benchmarks, including the datasets used in this study (Table 1), consist
214 of post-processed, 8-bit quantized images. These formats have undergone proprietary dynamic
215 range compression (log-transformation) and bit-depth reduction, which reduces the fine-grained
216 speckle structure. This reflects realistic clinical deployment conditions, where processed 8-bit
217 exports are the standard for hospital storage and telemedicine transmission.

218 While raw OCT acquisitions are influenced by multiple factors, including shot noise, ther-
219 mal noise, readout noise, speckle effects, reconstruction artifacts, and vendor-specific process-

220 ing, these components are no longer directly separable after image reconstruction and export.
 221 Therefore, the proposed multiplicative perturbation does not aim to physically simulate a specific
 222 OCT noise source. Instead, it provides a lightweight approximation of the residual acquisition-
 223 related variability observable in processed clinical OCT images, including variability arising from
 224 scanner hardware, reconstruction algorithms, image resolution, and vendor-specific preprocessing
 225 pipelines. Consequently, our noise injection strategy does not attempt to simulate raw interfero-
 226 metric speckle but rather models the residual speckle variability and scanner-specific noise texture
 227 remaining after clinical post-processing. Prior work has identified these residual texture differ-
 228 ences as a persistent source of cross-device domain shifts in vendor-exported OCT images.^{8,9} By
 229 injecting multiplicative intensity variations directly into the processed domain, we effectively tar-
 230 getted the device-dependent signal-to-noise signatures observable in deployed screening systems.
 231 This framing reflects the practical constraint that all existing large-scale OCT benchmarks operate
 232 exclusively in the post-acquisition image domain.

233 2.4.3 Noise-Aware Learning Formulation

234 Let $\mathbf{x} \in [0, 1]^{H \times W}$ where H and W denote the image height and width, respectively, denote the
 235 clean OCT B-scan. Noise-corrupted realizations are generated as

$$\mathbf{x}^{(n)} = \text{clip}(\mathbf{x} \odot (1 + \boldsymbol{\eta}^{(n)}), 0, 1), \quad \boldsymbol{\eta}^{(n)} \sim \mathcal{N}(0, \sigma^2), \quad (1)$$

236 where noise is injected in the post-export intensity domain prior to normalization to emulate the
 237 scanner-level variability observed in clinical OCT exports.

238 Per-image Z-normalization is defined as

$$\mathcal{Z}(\mathbf{x}) = \frac{\mathbf{x} - \mu(\mathbf{x})}{\sigma(\mathbf{x}) + \varepsilon}. \quad (2)$$

239 Each dynamic convolution layer maintains a bank of K expert kernels with effective filters

$$W_{\text{eff}}(\tilde{\mathbf{x}}) = \sum_{k=1}^K \alpha_k(\tilde{\mathbf{x}}) W_k, \quad \alpha(\tilde{\mathbf{x}}) = \text{Router}(\tilde{\mathbf{x}}). \quad (3)$$

240 Here, $\tilde{\mathbf{x}} = \mathcal{Z}(\mathbf{x})$ denotes the normalized OCT B-scan. Further, $\mu(\mathbf{x})$ and $\sigma(\mathbf{x})$ denote the mean
 241 and standard deviation of the pixel intensities in image \mathbf{x} , respectively, and ε is a small constant
 242 added for numerical stability to avoid division by zero during normalization. The index $n \in$
 243 $\{1, \dots, N\}$ denotes the stochastic noise realization, where N is the total number of realizations
 244 used during training. The function $f(\cdot)$ denotes the classifier mapping normalized OCT B-scans to
 245 class probabilities. Each dynamic convolution layer contains K expert kernels, where W_k denotes
 246 the k^{th} expert kernel and $\alpha_k(\cdot)$ denotes its corresponding routing weight. Finally, α represents the
 247 vector of routing weights produced by the router network.

248 The router is implemented as a lightweight attention module consisting of global average pool-
 249 ing, followed by two 1×1 convolutions with a ReLU nonlinearity and a softmax activation. This
 250 design yields per-image kernel mixture weights with negligible computational overhead, ensuring
 251 that dynamic routing does not compromise the inference efficiency.

252 2.4.4 Noise-Aware Objective

253 The training workflow is shown in Fig. 3. Classification is performed using the multi-class cross-
 254 entropy loss $\ell(\cdot, \cdot)$ with label smoothing to mitigate overconfidence under dataset shift. The stan-

255 dard loss is

$$\mathcal{L}_{\text{std}} = \ell(f(\mathcal{Z}(\mathbf{x})), y), \quad (4)$$

256 while the noise-aware loss is approximated as

$$\mathcal{L}_{\text{aug}} = \frac{1}{N} \sum_{n=1}^N \ell(f(\mathcal{Z}(\mathbf{x}^{(n)})), y). \quad (5)$$

257 The expectation is approximated via Monte-Carlo sampling over N independent realizations of
258 post-acquisition speckle variability per image.

259 The final training objective is

$$\mathcal{L} = \mathcal{L}_{\text{std}} + \beta \mathcal{L}_{\text{aug}}, \quad (6)$$

260 where β controls the strength of the noise-aware regularization.

261 Since multiplicative noise is injected prior to intensity clamping and per-image Z-normalization,
262 the expected value of each noisy realization does not equal the clean input. The clamping oper-
263 ation introduces distributional asymmetry, and subsequent Z-normalization further differentiates
264 the statistics of clean and noisy signals. The noise-aware objective therefore, exposes the routing
265 network to genuine input variance rather than mean-equivalent perturbations, encouraging kernel
266 mixtures that are stable across the scanner-dependent noise neighborhood of each B-scan.

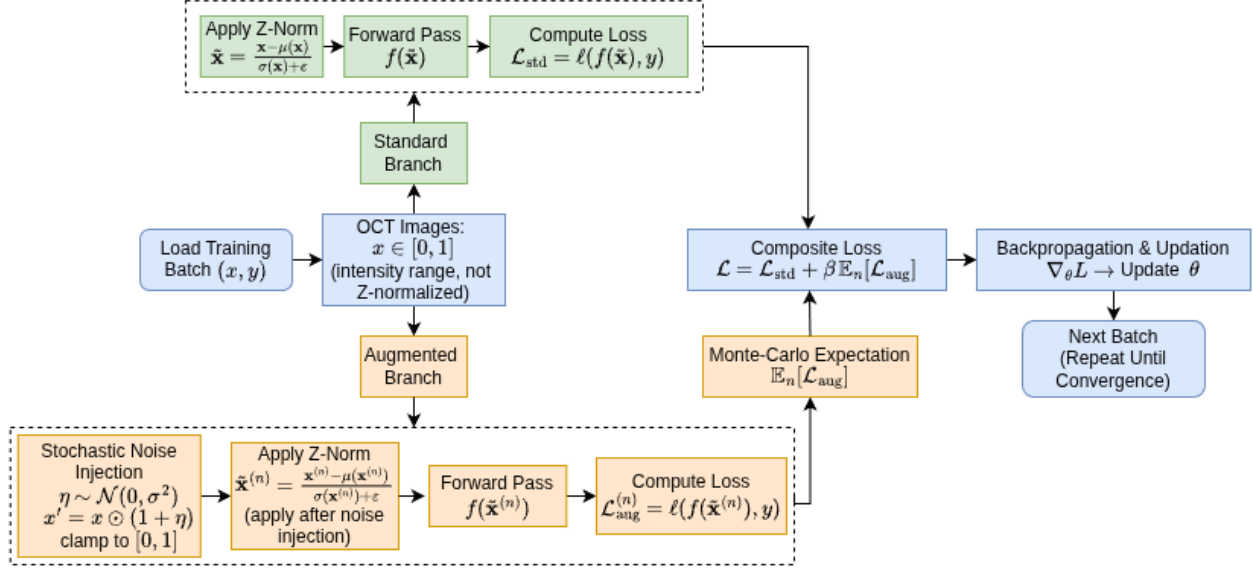


Fig 3 Training workflow of the proposed Noise-Aware Dynamic CNN. Each mini-batch is processed through a standard branch producing \mathcal{L}_{std} and a noise-augmented branch producing \mathcal{L}_{aug} , where multiplicative Gaussian noise is injected prior to per-image Z-normalization to simulate the post-acquisition scanner speckle variability. The Monte-Carlo estimate of the augmented loss is combined with the standard loss (Eq. 6) to update shared network parameters θ , encouraging the model to learn representations that are invariant to scanner-export speckle variability across devices.

267 2.4.5 Noise-Aware Dynamic Routing

268 For each clean image \mathbf{x} and each noisy realization $\mathbf{x}^{(n)}$, the routing network computes independent
 269 mixture weights

$$\alpha(\tilde{\mathbf{x}}) = \text{Router}(\tilde{\mathbf{x}}), \quad \alpha(\tilde{\mathbf{x}}^{(n)}) = \text{Router}(\tilde{\mathbf{x}}^{(n)}), \quad (7)$$

270 which define realization-specific effective kernels

$$W_{\text{eff}}(\tilde{\mathbf{x}}^{(n)}) = \sum_{k=1}^K \alpha_k(\tilde{\mathbf{x}}^{(n)}) W_k. \quad (8)$$

271 Gradients from the noise-aware loss propagate through the shared routing network, implicitly
 272 discouraging kernel-selection strategies that are oversensitive to scanner-specific perturbations.

273 Because both the actual and noise-corrupted realizations update the same routing parameters, the

274 model learns kernel mixtures that are more stable across stochastic perturbations. Unlike conven-
275 tional data augmentation, which perturbs inputs without explicitly influencing routing behavior, the
276 proposed objective indirectly promotes robustness in kernel selection through shared optimization
277 across multiple realizations of the same B-scan.

278 *2.4.6 Deployment Efficiency*

279 At inference time, noise sampling is disabled, and only a single clean forward pass is performed.
280 Consequently, NA-DyCNN retains computational complexity and latency identical to those of the
281 baseline Dynamic CNN while exhibiting substantially improved robustness to scanner and cohort
282 shift.

283 *2.5 Training Protocol*

284 All models were trained using the AdamW optimizer²⁶ with weight decay $\lambda = 10^{-4}$ and a batch
285 size of 64. The learning rate schedule employed a linear warm-up from 10^{-6} to 1×10^{-3} over
286 five epochs, followed by cosine annealing²⁷ to 10^{-6} over the remaining training period. Class
287 imbalance was addressed using inverse frequency weighting. Early stopping was applied based
288 on balanced validation accuracy with a patience of 10 epochs and a maximum training budget of
289 100 epochs. For each architecture, the hyperparameters were selected based on the performance
290 on the validation split of the source dataset and then fixed across all cross-dataset evaluations for
291 that model. The held-out test sets and all external datasets were never used during model selection
292 and were accessed only once for the final evaluation.

293

294 **Noise-Aware Hyperparameter Sensitivity.** The hyperparameters for the Noise-Aware Dynamic
 295 CNN, speckle variability scale (σ), regularization weight (β), and number of Monte Carlo realiza-
 296 tions (N), were tuned using only the source-dataset-balanced validation split, with no access to the
 297 target-domain data. For the D_2 -trained models, a grid search was performed over

$$\sigma \in \{0.05, 0.10, 0.15, 0.20\}, \quad \beta \in \{0.10, 0.30\}, \quad N \in \{1, 2, 4, 6, 8\},$$

298 with all other architectural and optimization hyperparameters fixed. The model was selected based
 299 on the source domain validation accuracy.

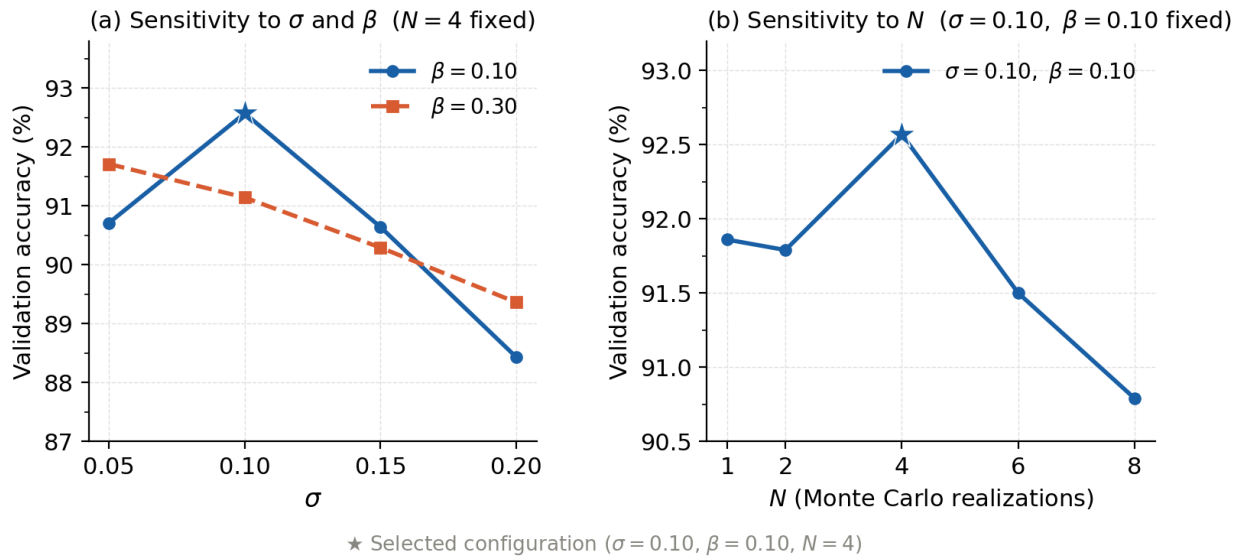


Fig 4 Hyperparameter sensitivity analysis on the OCT-C8 source-domain validation split. (a) Validation accuracy as a function of speckle scale σ for $\beta \in \{0.10, 0.30\}$ with $N = 4$ fixed. (b) Validation accuracy as a function of the number of Monte Carlo realizations N with $\sigma = 0.10$ and $\beta = 0.10$ fixed. The selected configuration ($\sigma = 0.10, \beta = 0.10, N = 4$) achieved the best validation performance among the explored configurations and was used for all cross-dataset evaluations.

300 Figure 4 shows the resulting sensitivity analysis for the D_2 source dataset. Moderate perturba-
 301 tion strengths yielded the strongest performance, with the validation accuracy peaking at $\sigma = 0.10$
 302 for $\beta = 0.10$. Increasing σ beyond this point led to a consistent decline in performance, suggesting

303 that excessively strong perturbations introduce unrealistic image corruption that degrades the diag-
304 nostically relevant retinal structure. Increasing the number of Monte Carlo realizations improved
305 the performance up to $N = 4$, after which the gains saturated and slightly declined despite the
306 increased computational cost. An identical tuning protocol was applied to the D_1 -trained models.
307 The initial coarse search identified the optimal region near $\sigma = 0.05$, and a subsequent finer local
308 search yielded the final configuration $\sigma = 0.06$ and $\beta = 0.12$. For both training datasets, the
309 validation performance was maximized at $N = 4$. Consequently, the final hyperparameters were
310 fixed to $(\sigma = 0.06, \beta = 0.12, N = 4)$ for D_1 and $(\sigma = 0.10, \beta = 0.10, N = 4)$ for D_2 for all
311 subsequent cross-dataset evaluations.

312 When the target scanner is entirely unknown at training time, a conservative default config-
313 uration ($\sigma = 0.10, \beta = 0.10, N = 4$) selected using source-domain validation data provides
314 a practical operating point, as the sensitivity analysis in Fig. 4 demonstrates stable performance
315 across moderate perturbation strengths. More broadly, the proposed framework is intended to
316 improve robustness to acquisition-related variability arising from scanner hardware, image resolu-
317 tion, reconstruction algorithms, and vendor-specific preprocessing pipelines rather than to model a
318 single scanner’s noise distribution exactly.

319 2.6 Evaluation Protocol

320 1. **Cross-Dataset Generalization (Strict Zero Shot).** Models are trained independently on D_1
321 or D_2 and evaluated on all remaining datasets under strict patient-level splits. Fine-tuning,
322 weight adaptation, or batch norm updating was not performed on any target domain. The
323 only permitted target domain access is the 10% labeled subset used exclusively for Mondrian
324 conformal calibration in the site-specific safety analysis. For $D_1/D_2 \rightarrow D_3$ evaluation, the

325 DME output logit was discarded since no corresponding labels exist in NEH. The remaining
 326 logits for CNV, Drusen, and Normal were renormalized prior to computing performance
 327 metrics.

328 **2. Performance Metrics.** We report performance metrics, such as Accuracy, Macro Sensitiv-
 329 ity, Macro Specificity, and recall-weighted F_2 score ($\beta = 2$, weighted average) to prioritize
 330 pathological recall in screening.

331 **3. Reliability and Uncertainty Metrics.** Model reliability was assessed using Cohen’s κ ²⁸ and
 332 Expected Calibration Error (ECE)²⁹ computed with 15 equal-mass bins.

333 The uncertainty quality was evaluated using the Error–Uncertainty metric (Err.Unc., %),
 334 which measures the fraction of misclassified samples lying above the 75th percentile of
 335 predictive entropy, indicating correct uncertainty attribution.

$$\text{Err.Unc.} = \frac{\sum_i \mathbf{1}(\hat{y}_i \neq y_i \wedge H(\mathbf{p}_i) > Q_{75})}{\sum_i \mathbf{1}(\hat{y}_i \neq y_i)} \times 100, \quad (9)$$

336 where $H(\mathbf{p}_i)$ denotes the predictive entropy of sample i and Q_{75} is the 75th percentile of en-
 337 tropy over the test set. Higher values indicate that model errors are preferentially associated
 338 with high uncertainty, reflecting safer clinical behavior.

339 **4. Mondrian Conformal Prediction.** To obtain finite-sample uncertainty guarantees under
 340 dataset shift, we apply class-conditional (Mondrian) conformal prediction.^{30,31} For each
 341 class y , nonconformity scores are defined as

$$\alpha_i^y = 1 - p_{i,y}, \quad (10)$$

342 computed on calibration samples whose true label is y . Class-specific thresholds q_y are
343 obtained from the $(1 - \epsilon)$ quantile of $\{\alpha_i^y\}$, yielding prediction sets

$$\mathcal{C}(\mathbf{p}) = \{y : 1 - p_y \leq q_y\}, \quad (11)$$

344 This guarantees a marginal class-conditional coverage of at least $1 - \epsilon = 95\%$.

345 Two deployment scenarios were evaluated: (a) Zero-shot, where source-domain validation
346 data are used for calibration and applied directly to target domains; and (b) Site-specific,
347 where 10% of the target domain patients are used for calibration with enforced $\min(n_y) \geq 5$
348 samples per class, reflecting realistic semi-supervised clinical adaptation.

349 We report the Conformal Coverage, Average Set Size, and Singleton Rate as performance
350 metrics.

351 **5. Statistical Validation.** The observed improvements were validated using hierarchical boot-
352 strap resampling³² with 2000 iterations, stratified at both the seed and patient levels. We
353 report 95% confidence intervals for differences in Accuracy, Weighted F_2 , and Selective
354 Risk at 80% coverage (Risk@80%). Statistical significance was declared when the confi-
355 dence interval excluded zero. All results are reported as mean \pm standard deviation across
356 three random seeds.

357 **Selective Risk at Fixed Coverage.** To evaluate performance under abstention, we compute Se-
358 lective Risk at 80% Coverage (Risk@80%).³³ Predictions were ranked by confidence (low pre-
359 dictive entropy), and the top 80% were retained. Risk@80% is the classification error for this

360 subset:

$$\text{Risk@80\%} = \frac{1}{|\mathcal{S}_{0.8}|} \sum_{i \in \mathcal{S}_{0.8}} \mathbf{1}(\hat{y}_i \neq y_i),$$

361 where $\mathcal{S}_{0.8}$ denotes the retained set. Lower values indicate safer, high-confidence predictions.

362 **3 Results**

363 We evaluated the clinical readiness of the proposed NA-DyCNN for cross-dataset deployment
364 through four lenses: (1) diagnostic accuracy across dataset cohorts, establishing baseline com-
365 petence; (2) formal safety via site-specific Mondrian conformal prediction; (3) reliability of un-
366 certainty estimates under dataset shift; and (4) computational efficiency, enabling real-time de-
367 ployment. This progression reflects deployment priorities; a model must diagnose accurately, fail
368 safely, and operate within realistic clinical hardware constraints.

Table 3 Cross-dataset diagnostic competence under domain shift (mean \pm std over three seeds). Bold indicates the best lightweight model. All datasets used are summarized in Table 1.

| Train \rightarrow Test | Model | Weighted F_2 Score (%) \uparrow | Accuracy (%) \uparrow | Sensitivity (%) \uparrow | Specificity (%) \uparrow |
|--------------------------|-----------------|-------------------------------------|---------------------------------|---------------------------------|---------------------------------|
| $D_1 \rightarrow D_3$ | Basic CNN | 69.65 \pm 2.4 | 71.27 \pm 1.9 | 68.91 \pm 2.3 | 85.57 \pm 0.7 |
| | Dynamic CNN | 66.78 \pm 8.9 | 69.33 \pm 8.3 | 65.91 \pm 6.6 | 84.62 \pm 3.0 |
| | NA-DyCNN | 72.54\pm2.1 | 74.22\pm2.5 | 70.66\pm1.1 | 86.55\pm0.6 |
| | MobileNetV3 | 71.88 \pm 3.2 | 72.84 \pm 3.7 | 70.08 \pm 1.5 | 85.29 \pm 1.0 |
| | ResNet18 (ref) | 76.35 \pm 3.5 | 77.75 \pm 3.2 | 73.57 \pm 2.9 | 87.41 \pm 1.3 |
| $D_2 \rightarrow D_3$ | Basic CNN | 62.15 \pm 5.5 | 62.90 \pm 5.2 | 64.54 \pm 4.5 | 82.15 \pm 1.9 |
| | Dynamic CNN | 71.06 \pm 4.0 | 71.80 \pm 4.4 | 72.13 \pm 1.7 | 85.90 \pm 1.8 |
| | NA-DyCNN | 74.90\pm2.3 | 75.49\pm2.3 | 74.52\pm1.2 | 87.34\pm0.6 |
| | MobileNetV3 | 62.92 \pm 6.7 | 63.93 \pm 6.4 | 62.65 \pm 3.4 | 81.65 \pm 1.8 |
| | ResNet18 (ref) | 78.34 \pm 1.9 | 78.52 \pm 1.8 | 78.08 \pm 2.3 | 88.74 \pm 0.9 |
| $D_2 \rightarrow D_1$ | Basic CNN | 87.01 \pm 2.2 | 86.97 \pm 2.2 | 86.90 \pm 2.3 | 95.83 \pm 0.7 |
| | Dynamic CNN | 89.98 \pm 1.9 | 89.95 \pm 2.0 | 90.75 \pm 2.2 | 96.83 \pm 0.7 |
| | NA-DyCNN | 92.89\pm0.4 | 92.87\pm0.4 | 93.02\pm0.6 | 97.71\pm0.1 |
| | MobileNetV3 | 90.75 \pm 1.3 | 90.72 \pm 1.3 | 90.60 \pm 1.2 | 97.09 \pm 0.4 |
| | ResNet18 (ref) | 94.41 \pm 0.9 | 94.40 \pm 0.9 | 94.61 \pm 0.9 | 98.21 \pm 0.3 |
| $D_1 \rightarrow D_2$ | Basic CNN | 91.18 \pm 0.3 | 91.21 \pm 0.3 | 91.21 \pm 0.3 | 97.07 \pm 0.1 |
| | Dynamic CNN | 92.05 \pm 0.6 | 92.09 \pm 0.6 | 92.09 \pm 0.6 | 97.36 \pm 0.2 |
| | NA-DyCNN | 92.29\pm0.6 | 92.33\pm0.6 | 92.33\pm0.6 | 97.44\pm0.2 |
| | MobileNetV3 | 91.96 \pm 1.2 | 91.98 \pm 1.2 | 91.98 \pm 1.2 | 97.33 \pm 0.4 |
| | ResNet18 (ref) | 94.33 \pm 0.5 | 94.34 \pm 0.5 | 94.34 \pm 0.5 | 98.11 \pm 0.2 |

Table 3 reports the diagnostic generalization across four cross-dataset transfer scenarios: $D_1 \rightarrow D_2$, $D_1 \rightarrow D_3$, $D_2 \rightarrow D_1$, and $D_2 \rightarrow D_3$, capturing a spectrum of acquisition-protocol, population, and annotation variability across hospital cohorts. Metrics were reported in a clinically prioritized order: weighted F_2 -score (recall-weighted), accuracy, sensitivity, and specificity.

Across all four transfer settings, NA-DyCNN consistently achieved the highest recall-weighted F_2 -scores among the lightweight models, indicating superior preservation of pathological sensitivity under dataset shift. Performance gains are most pronounced under $D_1 \rightarrow D_3$ and $D_2 \rightarrow D_3$, demonstrating that the proposed architecture maintains diagnostic competence even under substantial cross-site variability.

379 *3.2 Safety under Cross-Site Domain Shift*

380 Beyond diagnostic competence, we assessed clinical safety under cross-site domain shifts using
 381 Mondrian conformal prediction (MCP), a class-conditional uncertainty quantification framework
 382 that provides distribution-free coverage guarantees. We evaluated two deployment regimes: no-
 383 adaptation (zero-shot MCP using source-domain calibration only) and site-specific MCP with 10%
 384 labeled target-domain calibration, reflecting realistic clinical rollout conditions. The results are
 385 presented in Table 4.

Table 4 Mondrian conformal prediction safety under cross-site domain shift. Comparison of zero-shot (no target calibration) and site-specific (10% target calibration) models. Coverage is reported in %. Bold indicates the best lightweight model. All datasets used are summarized in Table 1.

| Shift | Model | No Adaptation (Zero-Shot) | | | Site-Specific (10% Calibration) | | |
|-----------------------|-----------------|----------------------------------|---------------------------------|----------------------------------|----------------------------------|---------------------------------|----------------------------------|
| | | Coverage (%) \uparrow | Set Size \downarrow | Singleton (%) \uparrow | Coverage (%) \uparrow | Set Size \downarrow | Singleton (%) \uparrow |
| $D_1 \rightarrow D_2$ | Basic CNN | 91.9 \pm 0.2 | 1.06 \pm 0.04 | 93.4 \pm 2.9 | 95.0 \pm 0.7 | 1.12 \pm 0.01 | 88.0 \pm 0.7 |
| | Dynamic CNN | 92.1 \pm 0.4 | 1.02 \pm 0.01 | 95.0 \pm 0.8 | 95.0 \pm 1.1 | 1.08 \pm 0.06 | 91.4 \pm 4.7 |
| | NA-DyCNN | 91.9 \pm 0.4 | 1.00\pm0.01 | 96.8\pm0.8 | 95.2\pm0.8 | 1.08\pm0.04 | 91.7\pm4.0 |
| | MobileNetV3 | 93.3\pm0.7 | 1.10 \pm 0.03 | 90.5 \pm 2.5 | 95.1 \pm 0.2 | 1.26 \pm 0.22 | 75.3 \pm 20.0 |
| | ResNet18 (ref) | 93.5 \pm 0.1 | 1.00 \pm 0.01 | 96.8 \pm 0.7 | 95.3 \pm 0.9 | 1.02 \pm 0.03 | 96.9 \pm 1.5 |
| $D_1 \rightarrow D_3$ | Basic CNN | 85.4 \pm 7.5 | 1.30 \pm 0.22 | 70.1 \pm 21.1 | 95.7\pm0.8 | 1.87 \pm 0.07 | 26.0 \pm 2.9 |
| | Dynamic CNN | 75.4 \pm 3.6 | 1.04\pm0.01 | 94.8\pm0.6 | 95.4 \pm 0.4 | 1.91 \pm 0.11 | 20.7 \pm 2.7 |
| | NA-DyCNN | 79.4 \pm 3.1 | 1.08 \pm 0.13 | 90.1 \pm 12.1 | 95.1 \pm 0.7 | 1.79\pm0.16 | 30.2\pm7.9 |
| | MobileNetV3 | 88.1\pm4.1 | 1.53 \pm 0.09 | 49.5 \pm 11.2 | 95.4 \pm 0.7 | 2.20 \pm 0.29 | 18.0 \pm 1.1 |
| | ResNet18 (ref) | 79.7 \pm 4.1 | 1.04 \pm 0.02 | 95.4 \pm 2.6 | 94.7 \pm 0.4 | 1.74 \pm 0.03 | 29.5 \pm 3.8 |
| $D_2 \rightarrow D_1$ | Basic CNN | 96.52\pm0.32 | 1.37 \pm 0.12 | 67.14 \pm 9.32 | 95.62 \pm 0.82 | 1.33 \pm 0.10 | 70.14 \pm 7.77 |
| | Dynamic CNN | 96.07 \pm 0.32 | 1.18 \pm 0.07 | 82.99 \pm 6.41 | 95.64\pm0.63 | 1.18 \pm 0.11 | 83.05 \pm 9.33 |
| | NA-DyCNN | 96.45 \pm 0.18 | 1.10\pm0.02 | 89.94\pm1.80 | 95.48 \pm 0.68 | 1.08\pm0.03 | 91.60\pm2.48 |
| | MobileNetV3 | 96.31 \pm 0.19 | 1.25 \pm 0.08 | 77.04 \pm 6.60 | 95.47 \pm 0.81 | 1.20 \pm 0.06 | 81.12 \pm 5.80 |
| | ResNet18 (ref) | 95.65 \pm 0.26 | 1.03 \pm 0.02 | 96.24 \pm 0.87 | 95.15 \pm 1.23 | 1.02 \pm 0.04 | 96.5 \pm 1.95 |
| $D_2 \rightarrow D_3$ | Basic CNN | 84.62\pm3.69 | 1.71 \pm 0.12 | 42.40 \pm 8.12 | 92.12 \pm 6.14 | 1.90 \pm 0.37 | 31.26 \pm 18.32 |
| | Dynamic CNN | 84.15 \pm 2.87 | 1.61 \pm 0.16 | 47.26 \pm 10.76 | 95.14\pm3.46 | 1.90 \pm 0.35 | 29.90 \pm 13.52 |
| | NA-DyCNN | 84.60 \pm 2.35 | 1.37\pm0.09 | 64.78\pm9.30 | 94.87 \pm 3.09 | 1.84\pm0.17 | 31.00\pm7.53 |
| | MobileNetV3 | 83.59 \pm 5.85 | 1.66 \pm 0.12 | 47.04 \pm 10.33 | 92.93 \pm 5.13 | 2.18 \pm 0.34 | 16.05 \pm 13.86 |
| | ResNet18 (ref) | 83.89 \pm 4.04 | 1.18 \pm 0.16 | 82.19 \pm 15.80 | 94.61 \pm 2.89 | 1.74 \pm 0.21 | 39.32 \pm 11.24 |

386 **Why Site-Specific Calibration is Necessary:** Table 4 shows that the effect of cross-dataset
 387 evaluation on conformal prediction strongly depends on the severity of the dataset mismatch. Un-
 388 der the $D_1 \rightarrow D_2$ transfer, all lightweight models maintained coverage relatively close to the nom-
 389 inal 95% target, even without target-domain calibration, indicating a relatively limited distribution

390 mismatch. In contrast, the more challenging transfers involving D_3 exhibit substantial zero-shot
391 under-coverage and inflated prediction sets, particularly for models without noise-aware training.
392 For example, under the severe $D_2 \rightarrow D_3$ shift, zero-shot deployment yielded only 84.6% cover-
393 age for NA-DyCNN, demonstrating the violation of exchangeability assumptions under a strong
394 domain shift. However, calibrating the MCP using only 10% of the labeled target-domain pa-
395 tients restored coverage to approximately 95% while maintaining compact prediction sets. This
396 demonstrates that minimal site-specific calibration is sufficient to recover the formal safety guar-
397 antees. Across the severe transfer settings, NA-DyCNN consistently achieved smaller prediction
398 sets and higher singleton rates than other lightweight baselines after calibration, indicating more
399 decisive and high-confidence predictions under cross-site variability. Empirically, larger calibra-
400 tion fractions further improved coverage and reduced the set size but were not pursued to avoid
401 increasing the clinical labeling burden, making the 10% setting a practical accuracy-cost trade-off
402 for deployment.

Table 5 Per-class Mondrian conformal prediction performance for the $D_1 \rightarrow D_2$ transfer scenario. Zero-shot corresponds to source-only calibration, whereas site-specific uses 10% target-domain calibration. Coverage and singleton values are reported in %. Datasets used are summarized in Table 1. Bold indicates the best lightweight model within each class.

| Class | Model | No Adaptation (Zero-Shot) | | | Site-Specific (10% Calibration) | | |
|--------|-----------------|--------------------------------|---------------------------------|--------------------------------|---------------------------------|---------------------------------|--------------------------------|
| | | Coverage (%) \uparrow | Set Size \downarrow | Singleton (%) \uparrow | Coverage (%) \uparrow | Set Size \downarrow | Singleton (%) \uparrow |
| CNV | Basic CNN | 92.1 \pm 0.7 | 1.07 \pm 0.03 | 92.0 \pm 2.5 | 93.8 \pm 1.5 | 1.08 \pm 0.02 | 91.6 \pm 1.0 |
| | Dynamic CNN | 93.0 \pm 1.0 | 1.04 \pm 0.01 | 94.6 \pm 0.9 | 94.1 \pm 1.4 | 1.04\pm0.05 | 94.9\pm3.0 |
| | NA-DyCNN | 93.1 \pm 0.6 | 1.03\pm0.02 | 95.8\pm1.5 | 94.3 \pm 1.7 | 1.06 \pm 0.04 | 94.3 \pm 3.9 |
| | MobileNetV3 | 94.9\pm1.1 | 1.08 \pm 0.05 | 92.0 \pm 5.3 | 95.6\pm1.4 | 1.57 \pm 0.53 | 47.9 \pm 47.4 |
| | ResNet18 (ref) | 94.5 \pm 1.2 | 1.02 \pm 0.01 | 97.4 \pm 0.2 | 94.7 \pm 2.6 | 1.01 \pm 0.03 | 97.3 \pm 0.7 |
| DME | Basic CNN | 96.9 \pm 0.1 | 1.02 \pm 0.01 | 97.9 \pm 1.1 | 95.3 \pm 1.1 | 1.10 \pm 0.01 | 89.7 \pm 1.4 |
| | Dynamic CNN | 97.0 \pm 0.1 | 1.00\pm0.00 | 98.1 \pm 0.3 | 96.0 \pm 0.6 | 1.05\pm0.03 | 94.1\pm2.4 |
| | NA-DyCNN | 96.9 \pm 0.3 | 1.00\pm0.00 | 98.7\pm0.3 | 96.6\pm0.3 | 1.07 \pm 0.03 | 92.6 \pm 2.5 |
| | MobileNetV3 | 97.6\pm0.4 | 1.02 \pm 0.01 | 97.8 \pm 0.6 | 95.3 \pm 0.1 | 1.11 \pm 0.07 | 89.5 \pm 6.9 |
| | ResNet18 (ref) | 97.2 \pm 0.3 | 1.00 \pm 0.00 | 98.9 \pm 0.2 | 96.0 \pm 0.5 | 1.00 \pm 0.02 | 98.6 \pm 0.9 |
| DRUSEN | Basic CNN | 96.8 \pm 0.5 | 1.06 \pm 0.02 | 94.0 \pm 2.0 | 95.8\pm1.4 | 1.16 \pm 0.01 | 84.3 \pm 1.0 |
| | Dynamic CNN | 97.1\pm0.2 | 1.04 \pm 0.02 | 95.5 \pm 1.8 | 95.3 \pm 2.4 | 1.14 \pm 0.11 | 85.9 \pm 10.0 |
| | NA-DyCNN | 96.6 \pm 0.2 | 1.02\pm0.01 | 97.3\pm1.0 | 95.7 \pm 2.0 | 1.11\pm0.08 | 89.2\pm7.3 |
| | MobileNetV3 | 96.9 \pm 0.8 | 1.07 \pm 0.03 | 93.4 \pm 3.1 | 94.9 \pm 1.4 | 1.23 \pm 0.16 | 78.8 \pm 13.4 |
| | ResNet18 (ref) | 97.4 \pm 0.4 | 1.02 \pm 0.01 | 97.5 \pm 0.5 | 95.9 \pm 0.7 | 1.05 \pm 0.03 | 95.1 \pm 3.1 |
| NORMAL | Basic CNN | 81.8 \pm 1.4 | 1.07 \pm 0.08 | 89.9 \pm 5.9 | 95.1\pm0.1 | 1.14 \pm 0.03 | 86.3 \pm 2.6 |
| | Dynamic CNN | 81.3 \pm 0.5 | 1.00 \pm 0.01 | 91.9 \pm 0.6 | 94.8 \pm 0.5 | 1.09 \pm 0.04 | 90.7 \pm 4.3 |
| | NA-DyCNN | 81.0 \pm 0.7 | 0.97\pm0.02 | 95.6\pm0.6 | 94.2 \pm 0.5 | 1.09\pm0.03 | 90.7\pm2.8 |
| | MobileNetV3 | 83.7\pm1.1 | 1.22 \pm 0.03 | 78.9 \pm 2.9 | 94.5 \pm 1.0 | 1.15 \pm 0.14 | 85.1 \pm 13.3 |
| | ResNet18 (ref) | 84.8 \pm 1.2 | 0.96 \pm 0.03 | 93.4 \pm 2.2 | 94.8 \pm 1.8 | 1.02 \pm 0.03 | 96.5 \pm 1.7 |

403 To examine class-level reliability, Table 5 reports per-class Mondrian conformal prediction
404 performance for the representative $D_1 \rightarrow D_2$ transfer setting. Consistent with the aggregate results,
405 NA-DyCNN generally produced the most compact prediction sets and among the highest singleton
406 rates across disease categories while maintaining coverage close to the nominal target. Similar
407 trends were observed for the remaining transfer settings and are provided in Supplementary Tables
408 S1–S3.

Table 6 Uncertainty calibration quality under cross-site domain shift (mean \pm std over three seeds). Bold text indicates the best performance among the lightweight models. Datasets used are summarized in Table 1. * indicates that the proposed approach outperforms ResNet18 on that metric.

| Shift | Model | Cohen’s κ \uparrow | ECE \downarrow | Err.Unc (%) \uparrow |
|-----------------------|-----------------|-----------------------------------|-----------------------------------|---------------------------------|
| $D_2 \rightarrow D_1$ | Basic CNN | 0.812 \pm 0.030 | 0.095 \pm 0.002 | 64.1 \pm 3.1 |
| | Dynamic CNN | 0.855 \pm 0.028 | 0.084 \pm 0.006 | 69.5 \pm 6.8 |
| | NA-DyCNN | 0.896\pm0.006 | 0.083 \pm 0.009 | 79.8\pm1.6* |
| | MobileNetV3 | 0.865 \pm 0.018 | 0.046\pm0.004 | 77.7 \pm 2.6 |
| | ResNet18 (Ref) | 0.918 \pm 0.013 | 0.058 \pm 0.008 | 79.7 \pm 3.3 |
| $D_2 \rightarrow D_3$ | Basic CNN | 0.436 \pm 0.067 | 0.056 \pm 0.038 | 37.6 \pm 4.0 |
| | Dynamic CNN | 0.558 \pm 0.056 | 0.035\pm0.007 | 47.4 \pm 5.2 |
| | NA-DyCNN | 0.608\pm0.029 | 0.042 \pm 0.002 | 48.0\pm1.4 |
| | MobileNetV3 | 0.434 \pm 0.074 | 0.158 \pm 0.048 | 36.3 \pm 3.8 |
| | ResNet18 (Ref) | 0.655 \pm 0.029 | 0.021 \pm 0.005 | 51.3 \pm 3.1 |
| $D_1 \rightarrow D_2$ | Basic CNN | 0.883 \pm 0.004 | 0.117 \pm 0.015 | 57.8 \pm 4.7 |
| | Dynamic CNN | 0.895 \pm 0.008 | 0.097 \pm 0.004 | 61.5\pm2.4 |
| | NA-DyCNN | 0.898\pm0.008 | 0.098 \pm 0.009 | 59.7 \pm 5.3* |
| | MobileNetV3 | 0.893 \pm 0.015 | 0.091\pm0.005 | 49.0 \pm 7.5 |
| | ResNet18 (Ref) | 0.925 \pm 0.006 | 0.106 \pm 0.009 | 57.7 \pm 4.3 |
| $D_1 \rightarrow D_3$ | Basic CNN | 0.540 \pm 0.028 | 0.063 \pm 0.032 | 41.3 \pm 1.8 |
| | Dynamic CNN | 0.512 \pm 0.110 | 0.070 \pm 0.036 | 39.6 \pm 4.3 |
| | NA-DyCNN | 0.580\pm0.030 | 0.050 \pm 0.014 | 44.4 \pm 2.7 |
| | MobileNetV3 | 0.557 \pm 0.046 | 0.037\pm0.005 | 45.3\pm4.6 |
| | ResNet18 (Ref) | 0.629 \pm 0.047 | 0.084 \pm 0.017 | 52.1 \pm 6.5 |

Cohen’s κ : inter-rater reliability beyond chance; **ECE**: expected calibration error (15 bins); **Err.Unc (%)**: fraction of misclassified samples exceeding the 75th percentile of predictive entropy, indicating correct uncertainty attribution.

410 As shown in Table 6, the uncertainty calibration behavior of the NA-DyCNN is shift-dependent.
 411 Under the more severe transfer settings involving D_3 , NA-DyCNN consistently achieved the high-
 412 est Cohen’s κ among the lightweight models together with the strongest uncertainty attribution
 413 performance (Err.Unc), indicating that prediction errors were more consistently associated with
 414 elevated predictive uncertainty. Under the milder $D_1 \rightarrow D_2$ and $D_2 \rightarrow D_1$ shifts, the lightweight
 415 models achieved broadly comparable calibration quality, with MobileNetV3 obtaining the lowest
 416 ECE. These results suggest that the proposed noise-aware routing mechanism provides the largest
 417 reliability benefits under substantial cross-site variability rather than mild domain shifts. More

418 generally, the results indicate that calibration quality alone does not fully capture deployment reli-
 419 ability under domain shift, thereby motivating the complementary use of conformal prediction to
 420 obtain formal uncertainty guarantees independent of probabilistic calibration quality.

421 3.3 Qualitative Error Analysis

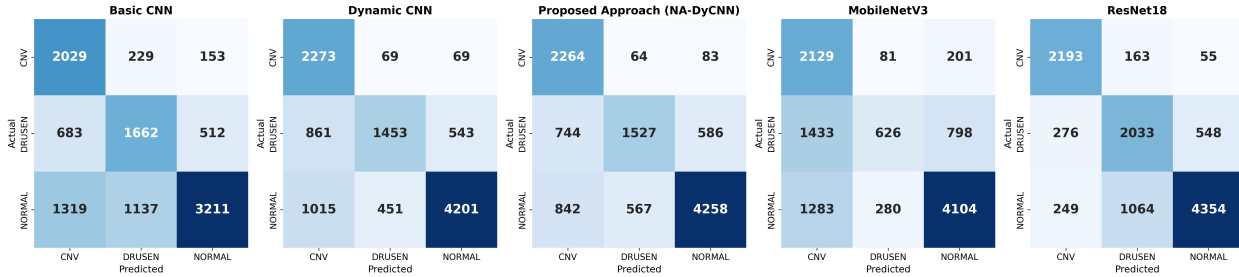


Fig 5 Confusion matrices for the cross-dataset shift ($D_2 \rightarrow D_3$). For each architecture, the model instance corresponding to the median macro-F1 score across three independent seeds is visualized.

422 Figure 5 reveals that Drusen remains the dominant source of diagnostic ambiguity under severe
 423 cross-dataset shift ($D_2 \rightarrow D_3$), consistent with its known morphological overlap with both early
 424 neovascular changes and healthy retinal appearance. All lightweight models exhibited bidirectional
 425 confusion between Drusen and CNV, as well as Drusen and Normal, highlighting the intrinsic
 426 difficulty of this class under a scanner-induced domain shift.

427 Importantly, NA-DyCNN does not merely improve aggregate performance but reshapes the
 428 error structure in clinically meaningful ways. Compared to the standard Dynamic CNN, NA-
 429 DyCNN reduces hazardous Drusen→CNV false positives from 861 to 744 cases (-13.6%) while
 430 simultaneously increasing Drusen recall (1453→1527 correct classifications), indicating improved
 431 sensitivity without sacrificing specificity. In parallel, NA-DyCNN also lowers Normal→CNV false
 432 alarms (1015→842), reducing unnecessary referrals in screening scenarios.

433 In contrast, MobileNetV3 exhibited a pronounced representation collapse under scanner shift,

434 misclassifying more than half of Drusen scans as CNV and achieving only 21.9% Drusen sensitiv-
 435 ity, showcasing the fragility of aggressive channel reduction strategies under the domain shift.

436 While ResNet18 achieves the highest Drusen recall (2033 correct classifications), it does so
 437 at nearly $28\times$ the parameter count of NA-DyCNN. Therefore, the proposed approach offers a
 438 favorable robustness–efficiency trade-off, substantially mitigating clinically hazardous confusions
 439 while retaining lightweight inference suitable for point-of-care OCT deployment.

440 3.4 Computational Efficiency for Clinical Deployment

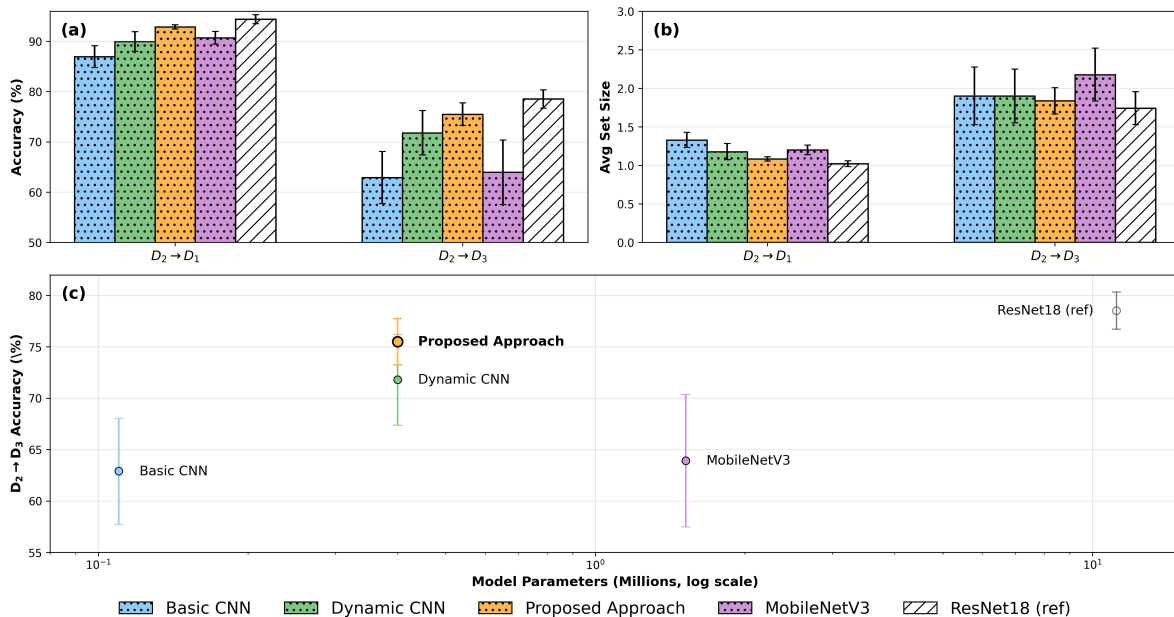


Fig 6 Clinical readiness metrics. (a) Diagnostic accuracy ($D_2 \rightarrow D_1$, $D_2 \rightarrow D_3$) with mean \pm std over 3 seeds. (b) Average Mondrian conformal set size under site-specific 10% calibration ($\epsilon = 0.05$), reflecting clinical deployment realism. Lower set size = more decisive prediction. (c) Accuracy–efficiency Pareto frontier: NA-DyCNN achieves an optimal balance among lightweight models. ResNet18 (hollow markers) shown as high-capacity reference only.

441 Figure 6(c) shows that NA-DyCNN lies on the accuracy–efficiency Pareto frontier among lightweight
 442 models, achieving 75.5% accuracy on the most severe cross-site shift ($D_2 \rightarrow D_3$) with only 0.40M
 443 parameters and 0.53 ms inference latency on an NVIDIA GB 10 GPU. This corresponds to a $27.9\times$

parameter reduction compared with ResNet18 (11.2M parameters), while retaining 95.6% of its diagnostic accuracy.

These results indicate that NA-DyCNN offers a favorable trade-off between predictive performance and computational cost, making it suitable for deployment in resource-constrained clinical environments, such as portable OCT workstations or edge-based screening systems, where memory footprint and latency are critical constraints.

3.5 Statistical Validation of Performance Gains

We assess robustness using hierarchical bootstrap resampling (2000 iterations, patient-stratified) across random seeds and the test subjects. We report differences in Selective Risk at 80% coverage (Risk@80%), recall-weighted F_2 score, and Accuracy between NA-DyCNN and each baseline. Confidence intervals excluding zero indicate statistical significance.

Table 7 Hierarchical bootstrap significance under scanner shift. The primary metric is the Selective Risk at 80% coverage (Risk@80%). A negative Δ Risk indicates safer, high-confidence predictions. All datasets used are summarized in Table 1.

| Shift | vs Baseline | Δ Risk@80% | ΔF_2 | Δ Accuracy |
|-----------------------|-------------|------------------------------|----------------------------|----------------------------|
| $D_1 \rightarrow D_2$ | Basic CNN | -0.75 [-2.13,0.33] | +1.10 [-0.12,1.92] | +1.11 [-0.10,1.93] |
| | Dynamic CNN | +0.06 [-1.28,1.03] | +0.24 [-0.10,0.60] | +0.23 [-0.10,0.58] |
| | MobileNetV3 | -1.42 [-2.74,0.79] | +0.33 [-2.07,2.03] | +0.35 [-2.02,2.02] |
| $D_2 \rightarrow D_1$ | Basic CNN | -4.21 [-6.28,-2.19] | +6.48 [4.00,9.19] | +5.90 [3.17,8.46] |
| | Dynamic CNN | -1.86 [-3.45,-0.26] | +2.70 [0.11,5.51] | +2.91 [0.82,5.22] |
| | MobileNetV3 | -1.72 [-2.76,-0.30] | +2.56 [0.47,4.23] | +2.14 [0.06,3.80] |
| $D_2 \rightarrow D_3$ | Basic CNN | -13.54 [-25.69,-3.39] | +11.48 [4.48,21.08] | +12.74 [5.20,22.10] |
| | Dynamic CNN | -3.50 [-4.96,-2.42] | +2.99 [1.56,4.10] | +3.70 [1.08,8.02] |
| | MobileNetV3 | -11.95 [-18.52,-6.57] | +12.49 [7.92,15.80] | +11.41 [6.38,16.64] |

Table 7 shows that the statistical significance of the NA-DyCNN improvements depends on the severity of the shift. For the mild $D_1 \rightarrow D_2$ transfer, improvements over Basic CNN, Dynamic CNN, and MobileNetV3 are small, and the corresponding confidence intervals overlap with zero,

458 indicating that the observed differences are not statistically significant. In contrast, for the more
 459 challenging $D_2 \rightarrow D_1$ and $D_2 \rightarrow D_3$ transfers, NA-DyCNN consistently achieves statistically sig-
 460 nificant reductions in selective diagnostic risk while simultaneously improving the recall-weighted
 461 F_2 score and classification accuracy.

462 The largest gains are observed under the $D_2 \rightarrow D_3$ transfer, where NA-DyCNN reduces the high-
 463 confidence error by 13.5 percentage points relative to the Basic CNN, 11.95 percentage points
 464 relative to MobileNetV3, and 3.5 percentage points relative to Dynamic CNN. These results sug-
 465 gest that the benefits of the proposed noise-aware objective become most pronounced under larger
 466 shifts, where robustness and calibration are particularly important for reliable deployment.

467 3.6 Mechanisms of Robustness: Ablation Study

468 To identify the individual contributions of each architectural and training component, we con-
 469 ducted an ablation study on the most challenging transfer ($D_2 \rightarrow D_3$). Table 8 reports the per-
 470 formance (mean \pm std over three seeds) as successive components are added, using site-specific
 471 Mondrian calibration with 10% target-domain patients. Variant (3) is included only as an inter-
 472 mediate ablation configuration to isolate the contribution of post-acquisition speckle augmentation
 473 independently from the expectation-based noise-aware objective.

Table 8 Ablation study on severe cross-site shift ($D_2 \rightarrow D_3$). Each component was added cumulatively to quantify its contribution to robustness. Descriptions of the datasets are summarized in Table 1.

| Variant | DynConv | Post-Acq. Noise Aug | Noise-Aware | Weighted F_2 (%) \uparrow | Accuracy (%) \uparrow | Cohen’s κ \uparrow | Set Size \downarrow |
|--------------------------------|---------|---------------------|-------------|---------------------------------|---------------------------------|-----------------------------------|---------------------------------|
| (1) Basic CNN | × | × | × | 62.15 \pm 5.5 | 62.90 \pm 5.2 | 0.434 \pm 0.067 | 1.9 \pm 0.37 |
| (2) Dynamic CNN | ✓ | × | × | 71.06 \pm 4.0 | 71.80 \pm 4.4 | 0.558 \pm 0.056 | 1.89 \pm 0.35 |
| (3) Dynamic CNN with Noise Aug | ✓ | ✓ | × | 72.85 \pm 1.6 | 72.85 \pm 1.8 | 0.572 \pm 0.025 | 1.95 \pm 0.11 |
| (4) Full NA-DyCNN | ✓ | ✓ | ✓ | 74.90\pm2.3 | 75.49\pm2.3 | 0.608\pm0.029 | 1.84\pm0.17 |

Values are mean \pm std over three random seeds under the $D_2 \rightarrow D_3$ shift with site-specific Mondrian conformal calibration (10% target-domain patients). **DynConv**: Dynamic convolution kernels (K=4, attention routing). **Post-Acq. Noise Aug**: Augmentation modeling residual scanner-export noise and texture variability in post-processed OCT images. **Noise-Aware**: Expectation-based loss (Eq. 6) with $N = 4$ realizations.

474 Table 8 shows that Dynamic convolution yields the largest absolute improvement, increasing
475 weighted F_2 by approximately +8.91 percentage points and accuracy by +8.9 percentage points
476 over the Basic CNN baseline, together with a substantial increase in Cohen’s κ (from 0.43 to 0.56).
477 Adding speckle-specific augmentation provides an additional boost of +1.8 percentage points in
478 weighted F_2 and modest gains in accuracy and κ , although with a slightly higher set size. The full
479 noise-aware objective produces the strongest overall gains, with a further +2.1 percentage point
480 improvement in weighted F_2 , +2.6 percentage points in accuracy, a +0.036 increase in κ , and the
481 lowest set size (1.84), indicating more decisive and reliable high-confidence predictions under the
482 cross-site domain shift.

483 These results indicate that robustness emerges from the combined effect of dynamic kernel
484 adaptation, OCT-specific speckle augmentation, and expectation-based optimization over speckle
485 perturbations rather than from any single component in isolation.

486

487 3.7 Analysis of Dataset Information in Routing Coefficients

488 To directly evaluate whether the proposed noise-aware objective reduces the dataset-specific in-
489 formation encoded in the routing decisions, we analyzed the routing coefficients (α) from the
490 final dynamic convolution block for the Dynamic CNN and NA-DyCNN. The D_1 -trained model
491 was utilized for this analysis because both D_1 and D_3 were acquired using Heidelberg Spectralis
492 scanners; therefore, distinguishing between these two datasets would provide limited insight into
493 scanner-related variability. Instead, the routing coefficients from D_2 and D_3 were evaluated, where
494 D_2 (OCT-C8) contained images acquired from multiple scanner vendors and D_3 consisted of im-
495 ages acquired using a Heidelberg Spectralis scanner. Consequently, distinguishing D_2 from D_3

496 provides a more meaningful assessment of whether the routing coefficients retain dataset- and
497 scanner-related information.

498 A logistic regression classifier was trained to predict dataset-of-origin (D_2 vs. D_3) using only the
499 routing coefficients α . Five-fold stratified cross-validation was performed, and the results were
500 averaged across three independent training seeds. Dynamic CNN achieved a mean dataset classifi-
501 cation accuracy of 0.756 ± 0.007 , whereas NA-DyCNN achieved 0.733 ± 0.008 . Because a lower
502 classification accuracy indicates that the dataset-of-origin is more difficult to infer from the rout-
503 ing coefficients, these results suggest that the proposed noise-aware objective reduces the dataset-
504 specific information encoded within the routing space. This observation is consistent with the
505 hypothesis that NA-DyCNN promotes more dataset-invariant routing behavior, which contributes
506 to the improved cross-dataset generalization results reported in Table 8.

507 **4 Discussion**

508 This study addresses a key barrier to the deployment of deep learning for retinal OCT screen-
509 ing: the limited generalization of lightweight models across heterogeneous scanners and clinical
510 sites. By incorporating post-acquisition residual speckle variability into a dynamic convolution
511 framework, NA-DyCNN improves cross-dataset robustness without increasing the inference cost
512 or requiring any form of target-domain fine-tuning. From an imaging perspective, this approach
513 explicitly accounts for the variability of OCT speckle statistics across acquisition systems. Speckle
514 patterns arise from coherent interference of backscattered light and are influenced by optical sys-
515 tem parameters and tissue microstructure; therefore, variations in scanner design and reconstruc-
516 tion pipelines naturally lead to differences in observed speckle distributions.³ However, scanner
517 variability is not the sole source of distribution shift in real-world OCT deployment. Differences in

518 patient demographics, acquisition protocols, export pipelines, and dataset curation strategies can
519 substantially alter the target distribution. This is particularly evident in the $D_1 \rightarrow D_3$ transfer set-
520 ting, where both datasets use Heidelberg Spectralis devices but still exhibit a substantial cross-site
521 distribution shift.

522 Across all external validation scenarios, NA-DyCNN achieved the best performance among
523 lightweight models under strict zero-shot transfer, where no model parameters or normalization
524 statistics were adapted to the target domain. Statistically significant gains are observed under the
525 most severe shifts ($D_2 \rightarrow D_3$ and $D_2 \rightarrow D_1$). Ablation analysis confirmed that dynamic convolu-
526 tion provided the dominant performance improvement, whereas the expectation-based noise-aware
527 objective further stabilized kernel routing against scanner-export speckle variability. These find-
528 ings suggest that explicitly modeling speckle perturbations during training helps prevent the model
529 from relying on scanner-specific texture signatures that do not generalize across acquisition sys-
530 tems.

531 Beyond diagnostic accuracy, NA-DyCNN improves clinical reliability, yielding higher Cohen’s
532 κ and more compact conformal prediction sets after minimal site-specific calibration. Although the
533 expected calibration error remains elevated under severe dataset shifts, this limitation is mitigated
534 by the use of conformal prediction, which provides distribution-free safety guarantees independent
535 of the probability calibration quality.

536 The under-coverage observed in the NEH cohort highlights the limitations of generalizability
537 induced by the combined scanner, population, and protocol shift, highlighting the inherent diffi-
538 culty of uncertainty quantification in real-world deployment. These findings are consistent with
539 prior observations that OCT speckle patterns encode acquisition-dependent information related to
540 both imaging physics and tissue scattering properties.³ As a result, models trained on a single

541 dataset may inadvertently learn scanner-specific speckle statistics rather than pathology-relevant
542 structures.

543 Recent lightweight architectures, including modern mobile CNNs (e.g., MobileNetV4), lightweight
544 transformer models, and hybrid CNN-transformer architectures such as MobileViT, EfficientViT,
545 and FastViT, have demonstrated strong performance across various computer vision tasks. How-
546 ever, the primary contribution of this study lies in the proposed noise-aware dynamic routing frame-
547 work rather than the backbone architecture itself. A deliberate design choice in this study was
548 to train all models from scratch rather than relying on large-scale pre-trained models. Although
549 transfer learning often improves performance, pretrained representations are typically learned from
550 natural image datasets and may introduce biases that are not directly relevant to retinal OCT imag-
551 ing. Moreover, transfer learning remains susceptible to domain shifts, making it difficult to dis-
552 entangle the effects of pre-training from those of the proposed robustness framework. The pro-
553 posed approach is architecture-agnostic and can be viewed as a general strategy for encouraging
554 scanner-invariant feature utilization via noise-aware routing. Although this study focuses on mod-
555 eling dataset-dependent speckle variability, other acquisition-related factors, such as resolution
556 differences, reconstruction pipelines, contrast variations, and device-specific artifacts, may also
557 contribute to domain shift. Future work could therefore investigate alternative augmentation and
558 perturbation strategies within the same framework, as well as integration with modern lightweight
559 CNN, transformers, and hybrid architectures.

560 Future work will also include validation on raw 12/16-bit OCT data to better disentangle
561 physical speckle formation from vendor post-processing artifacts, development of scanner-aware
562 calibration strategies, and extension to volumetric OCT and additional vendor platforms. More
563 broadly, the results support the emerging perspective that speckle should not be treated purely as

564 a nuisance artifact but as an imaging characteristic whose statistical variability must be considered
565 when developing robust OCT analysis algorithms.

566 Despite these limitations, NA-DyCNN retains an identical inference cost to standard Dynamic
567 CNNs while being nearly $28\times$ smaller than ResNet18, making it well-suited for real-time deploy-
568 ment in portable and resource-constrained OCT screening environments.

569 **5 Conclusion**

570 We presented NA-DyCNN, a noise-aware dynamic convolutional framework for retinal OCT clas-
571 sification that integrates scanner-export residual speckle variability into the training objective within
572 the clinical image-processing domain. By optimizing over multiple noise-corrupted realizations
573 that approximate variations in OCT speckle statistics across acquisition systems, the model im-
574 proves robustness to dataset shifts while preserving lightweight inference efficiency.

575 A comprehensive evaluation across three heterogeneous OCT cohorts demonstrated that NA-
576 DyCNN consistently outperformed existing lightweight baselines and approached the performance
577 of deeper residual networks at a fraction of the computational cost. Beyond accuracy, the proposed
578 framework improves diagnostic reliability and conformal prediction efficiency under dataset shifts,
579 supporting safer clinical deployment.

580 These findings suggest that explicitly accounting for speckle variability can improve the gen-
581 eralizability of deep learning models across scanners and clinical sites. More broadly, incorporat-
582 ing clinically observable imaging variability into adaptive convolutional architectures provides a
583 promising pathway toward robust and deployable OCT screening systems in real-world healthcare
584 environments.

585 *Disclosures*

586 The authors declare that there are no financial interests, commercial affiliations, or other potential
587 conflicts of interest that could have influenced the objectivity of this research or the writing of this
588 paper.

589 *Code and Data Availability*

590 All datasets used in this study are publicly available from their respective sources. The UCSD reti-
591 nal OCT dataset is available at: <https://data.mendeley.com/datasets/rsbjbr9sj/>
592 2. The NEH retinal OCT dataset can be accessed at: [https://data.mendeley.com/](https://data.mendeley.com/datasets/8kt969dhx6/2)
593 [datasets/8kt969dhx6/2](https://data.mendeley.com/datasets/8kt969dhx6/2). The OCT-C8 dataset is available via Kaggle at: [https://](https://www.kaggle.com/datasets/obulisainaren/retinal-oct-c8)
594 www.kaggle.com/datasets/obulisainaren/retinal-oct-c8. All datasets were
595 used in accordance with their respective usage policies and licenses.

596 *References*

- 597 1 D. Huang, E. A. Swanson, C. P. Lin, *et al.*, “Optical coherence tomography,” *science*
598 **254**(5035), 1178–1181 (1991).
- 599 2 W. Drexler and J. G. Fujimoto, *Optical coherence tomography: technology and applications*,
600 Springer (2008).
- 601 3 V. B. Silva, D. Andrade De Jesus, S. Klein, *et al.*, “Signal-carrying speckle in optical coher-
602 ence tomography: a methodological review on biomedical applications,” *Journal of biomed-*
603 *ical optics* **27**(3), 030901–030901 (2022).
- 604 4 F. Dhaoui and A. Zrelli, “Retinal diseases classification system using oct images combined

- 605 with cnn models,” in *2023 International Symposium on Networks, Computers and Commu-*
606 *nications (ISNCC)*, 1–6, IEEE (2023).
- 607 5 O. F. Aydın, M. S. Nazlı, F. B. Tek, *et al.*, “Retinal disease classification using optical coher-
- 608 ence tomography angiography images,” in *2024 9th International Conference on Computer*
609 *Science and Engineering (UBMK)*, 884–889, IEEE (2024).
- 610 6 S. Akça, Z. Garip, E. Ekinci, *et al.*, “Automated classification of choroidal neovascularization,
- 611 diabetic macular edema, and drusen from retinal oct images using vision transformers: a
- 612 comparative study,” *Lasers in Medical Science* **39**(1), 140 (2024).
- 613 7 D. Kermany, K. Zhang, and M. Goldbaum, “Large dataset of labeled optical coherence to-
- 614 mography (oct) and chest x-ray images.” Mendeley Data, Version 3 (2018).
- 615 8 J. Wang, Y. Chen, W. Li, *et al.*, “Domain adaptation model for retinopathy detection from
- 616 cross-domain oct images,” in *Proceedings of the Third Conference on Medical Imaging with*
617 *Deep Learning*, T. Arbel, I. Ben Ayed, M. de Bruijne, *et al.*, Eds., *Proceedings of Machine*
618 *Learning Research* **121**, 795–810, PMLR (2020).
- 619 9 F.-E. Jannat, S. Gholami, *et al.*, “Oct-selfnet: a self-supervised framework with multi-source
- 620 datasets for generalized retinal disease detection,” *Front. Big Data* , 1609124 (2025).
- 621 10 A. F. Fercher, W. Drexler, C. K. Hitzenberger, *et al.*, “Optical coherence tomography-
- 622 principles and applications,” *Reports on progress in physics* **66**(2), 239 (2003).
- 623 11 N. A. Kande, R. Dakhane, A. Dukkipati, *et al.*, “Siamesegan: A generative model for de-
- 624 noising of spectral domain optical coherence tomography images,” *IEEE Transactions on*
625 *Medical Imaging* **40**(1), 180–192 (2023).

- 626 12 Z. Baharlouei, H. Rabbani, and G. Plonka, “Wavelet scattering transform application in clas-
627 sification of retinal abnormalities using oct images,” *Scientific reports* **13**(1), 19013 (2023).
- 628 13 N. Paluru, H. Ravishankar, S. Hegde, *et al.*, “Self distillation for improving the generalizabil-
629 ity of retinal disease diagnosis using optical coherence tomography images,” *IEEE Journal*
630 *of Selected Topics in Quantum Electronics* **29**(4), 7200812 (2023).
- 631 14 A. Howard, M. Sandler, B. Chen, *et al.*, “Searching for mobilenetv3,” in *2019 IEEE/CVF*
632 *International Conference on Computer Vision (ICCV)*, 1314–1324 (2019).
- 633 15 K. He, X. Zhang, S. Ren, *et al.*, “Deep residual learning for image recognition,” in *2016 IEEE*
634 *Conference on Computer Vision and Pattern Recognition (CVPR)*, 770–778 (2016).
- 635 16 S. Agarwal, A. K. Dohare, P. Saxena, *et al.*, “Hdl-aco hybrid deep learning and ant colony op-
636 timization for ocular optical coherence tomography image classification,” *Scientific Reports*
637 **15**(1), 5888 (2025).
- 638 17 Y. Chen, X. Dai, M. Liu, *et al.*, “Dynamic convolution: Attention over convolution kernels,”
639 in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,
640 11030–11039 (2020).
- 641 18 O. S. Naren, “Retinal oct image classification - c8 dataset.” Kaggle Dataset (2021).
- 642 19 S. Sotoudeh-Paima, F. Hajizadeh, and H. Soltanian-Zadeh, “Labeled retinal optical coherence
643 tomography dataset for classification of normal, drusen, and cnv cases.” Mendeley Data,
644 Version 2 (2023).
- 645 20 S. H. Park and K. Han, “Methodologic guide for evaluating clinical performance and effect
646 of artificial intelligence technology for medical diagnosis and prediction,” *Radiology* **286**(3),
647 800–809 (2018).

- 648 21 C. E. Duchon, “Lanczos filtering in one and two dimensions,” *Journal of Applied Meteorol-*
649 *ogy (1962-1982)*, 1016–1022 (1979).
- 650 22 A. Paszke, S. Gross, F. Massa, *et al.*, “Pytorch: An imperative style, high-performance deep
651 learning library,” *Advances in neural information processing systems* **32** (2019).
- 652 23 L. Terry, “Topcon 3D-OCT 1000 repeatability dataset.tif,” (2016).
- 653 24 L. Terry, “Cirrus HD-OCT repeatability dataset.tif,” (2016).
- 654 25 M. Kulyabin, A. Zhdanov, A. Nikiforova, *et al.*, “Octl: Optical coherence tomography
655 dataset for image-based deep learning methods,” *Scientific data* **11**(1), 365 (2024).
- 656 26 I. Loshchilov and F. Hutter, “Decoupled weight decay regularization,” in *International Con-*
657 *ference on Learning Representations*, (2019).
- 658 27 I. Loshchilov and F. Hutter, “Sgdr: Stochastic gradient descent with warm restarts,” (2016).
- 659 28 J. Cohen, “A coefficient of agreement for nominal scales,” *Educational and psychological*
660 *measurement* **20**(1), 37–46 (1960).
- 661 29 C. Guo, G. Pleiss, Y. Sun, *et al.*, “On calibration of modern neural networks,” in *International*
662 *conference on machine learning*, 1321–1330, PMLR (2017).
- 663 30 H. Boström, U. Johansson, and T. Löfström, “Mondrian conformal predictive distributions,”
664 in *Proceedings of the Tenth Symposium on Conformal and Probabilistic Prediction and Ap-*
665 *plications*, L. Carlsson, Z. Luo, G. Cherubin, *et al.*, Eds., *Proceedings of Machine Learning*
666 *Research* **152**, 24–38, PMLR (2021).
- 667 31 A. N. Angelopoulos and S. Bates, “A gentle introduction to conformal prediction and
668 distribution-free uncertainty quantification,” *arXiv preprint arXiv:2107.07511* (2021).

669 32 B. Efron and R. J. Tibshirani, *An introduction to the bootstrap*, Chapman and Hall/CRC
670 (1994).

671 33 Y. Geifman and R. El-Yaniv, “Selective classification for deep neural networks,” *Advances in*
672 *neural information processing systems* **30** (2017).

673 **Deeksha Chutani** received her B.E. in Computer Engineering from Thapar Institute of Engineer-
674 ing and Technology, Patiala, and her B.S. in Data Science and Applications from the Indian Insti-
675 tute of Technology Madras in 2024. She is currently pursuing an M.Tech. in Computational and
676 Data Science at the Indian Institute of Science, Bengaluru. Her research interests lie in artificial
677 intelligence for the social good, with a particular focus on healthcare and medical imaging.

678 **Phaneendra K. Yalavarthy** holds a PhD from Dartmouth College and leads TANUH, the AI-
679 Center of Excellence in Healthcare, and is also a professor of Computational and Data Sciences
680 at IISc Bangalore. Awards include INAE Young Engineer (2013) and S. Ramachandran-National
681 Bioscience Award (2020). His research interests include computational/deep learning in medical
682 imaging, physiological signals, and nondestructive imaging. A senior IEEE/OSA/SPIE member,
683 he also serves as an associate editor of IEEE Transactions on Medical Imaging and PLOS Digital
684 Health.

685 **List of Tables**

686 1 [OCT dataset characteristics and usage protocols after quality control \(QC\).](#)

- 687 2 Computational complexity and inference efficiency for 256×256 grayscale input
688 on NVIDIA GB 10 with batch size 1. Latency measurements represent the mean \pm
689 std over 1000 forward passes. The relative training cost is normalized to Dynamic
690 CNN ($1.0\times$), measured on the D_2 training split (batch size 64). Bold entries denote
691 the proposed method.
- 692 3 Cross-dataset diagnostic competence under domain shift (mean \pm std over three
693 seeds). Bold indicates the best lightweight model. All datasets used are summa-
694 rized in Table 1.
- 695 4 Mondrian conformal prediction safety under cross-site domain shift. Comparison
696 of zero-shot (no target calibration) and site-specific (10% target calibration) mod-
697 els. Coverage is reported in %. Bold indicates the best lightweight model. All
698 datasets used are summarized in Table 1.
- 699 5 Per-class Mondrian conformal prediction performance for the $D_1 \rightarrow D_2$ transfer
700 scenario. Zero-shot corresponds to source-only calibration, whereas site-specific
701 uses 10% target-domain calibration. Coverage and singleton values are reported in
702 %. Datasets used are summarized in Table 1. Bold indicates the best lightweight
703 model within each class.
- 704 6 Uncertainty calibration quality under cross-site domain shift (mean \pm std over
705 three seeds). Bold text indicates the best performance among the lightweight mod-
706 els. Datasets used are summarized in Table 1. * indicates that the proposed ap-
707 proach outperforms ResNet18 on that metric.

- 708 7 Hierarchical bootstrap significance under scanner shift. The primary metric is the
709 Selective Risk at 80% coverage (Risk@80%). A negative Δ Risk indicates safer,
710 high-confidence predictions. All datasets used are summarized in Table 1.
- 711 8 Ablation study on severe cross-site shift ($D_2 \rightarrow D_3$). Each component was added
712 cumulatively to quantify its contribution to robustness. Descriptions of the datasets
713 are summarized in Table 1.

714 List of Figures

- 715 1 Visualization of OCT appearance variability, together with sample synthetic per-
716 turbations used during noise-aware training. (a) Representative OCT B-scans ob-
717 tained from publicly available sources acquired using Topcon,²³ Zeiss Cirrus,²⁴ and
718 Optovue Avanti RTVue XR²⁵ systems, illustrating the diversity of image appear-
719 ance, contrast characteristics, and retinal layer representation across different ven-
720 dors. (b) Representative OCT B-scans from the three datasets used in this study:
721 D_1 (UCSD; Heidelberg Spectralis), D_2 (OCT-C8; mixed vendors), and D_3 (NEH;
722 Heidelberg Spectralis), illustrating scanner-dependent variations in the speckle tex-
723 ture and reflectivity statistics. (c) Noise-perturbed realizations of the same OCT
724 B-scan generated using multiplicative Gaussian perturbations ($\sigma = 0.05, 0.10,$ and
725 0.15), illustrating increasing levels of synthetic speckle variability. These values
726 are shown for visualization purpose and do not correspond exactly to the hyperpa-
727 rameters used during training.

728 2 Structural comparison of convolutional architectures for OCT classification. (a)
729 Basic CNN applies per-image z -normalization, followed by static convolution with
730 shared kernels across all inputs. (b) Dynamic CNN employs a router network to
731 generate image-dependent kernel mixtures for dynamic convolution. (c) Noise-
732 Aware Dynamic CNN (Proposed) introduces expectation-based regularization by
733 training on multiple stochastic realizations of scanner-export speckle variability
734 $\mathbf{x}^{(n)}$ of each input, while using the same dynamic routing mechanism for both
735 clean and noisy samples. During inference, the noise branch is disabled, yielding
736 a computational cost identical to that of the standard Dynamic CNN. Here, $\ell(\cdot, \cdot)$
737 denotes the multi-class cross-entropy loss.

738 3 Training workflow of the proposed Noise-Aware Dynamic CNN. Each mini-batch
739 is processed through a standard branch producing \mathcal{L}_{std} and a noise-augmented
740 branch producing \mathcal{L}_{aug} , where multiplicative Gaussian noise is injected prior to per-
741 image Z -normalization to simulate the post-acquisition scanner speckle variability.
742 The Monte-Carlo estimate of the augmented loss is combined with the standard
743 loss (Eq. 6) to update shared network parameters θ , encouraging the model to learn
744 representations that are invariant to scanner-export speckle variability across de-
745 vices.

- 746 4 Hyperparameter sensitivity analysis on the OCT-C8 source-domain validation split.
- 747 (a) Validation accuracy as a function of speckle scale σ for $\beta \in \{0.10, 0.30\}$ with
- 748 $N = 4$ fixed. (b) Validation accuracy as a function of the number of Monte Carlo
- 749 realizations N with $\sigma = 0.10$ and $\beta = 0.10$ fixed. The selected configuration
- 750 ($\sigma = 0.10, \beta = 0.10, N = 4$) achieved the best validation performance among the
- 751 explored configurations and was used for all cross-dataset evaluations.
- 752 5 Confusion matrices for the cross-dataset shift ($D_2 \rightarrow D_3$). For each architecture,
- 753 the model instance corresponding to the median macro-F1 score across three inde-
- 754 pendent seeds is visualized.
- 755 6 Clinical readiness metrics. (a) Diagnostic accuracy ($D_2 \rightarrow D_1, D_2 \rightarrow D_3$) with
- 756 mean \pm std over 3 seeds. (b) Average Mondrian conformal set size under site-
- 757 specific 10% calibration ($\epsilon = 0.05$), reflecting clinical deployment realism. Lower
- 758 set size = more decisive prediction. (c) Accuracy-efficiency Pareto frontier: NA-
- 759 DyCNN achieves an optimal balance among lightweight models. ResNet18 (hol-
- 760 low markers) shown as high-capacity reference only.