

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

Deeksha Chutani and Phaneendra K. Yalavarthy

1 Class-Conditional Mondrian Conformal Prediction Analysis

To complement the aggregate Mondrian conformal prediction (MCP) results reported in the main manuscript, this section provides a detailed class-level analysis for the remaining cross-dataset transfer settings. For each transfer scenario, we report class-wise coverage, average prediction set size, and singleton rate under both zero-shot deployment and site-specific calibration using 10% labeled target-domain data. These results provide additional insight into how uncertainty estimates behave across individual retinal disease categories under dataset shift.

Table S1. Per-class Mondrian conformal prediction performance for the $D_1 \rightarrow D_3$ transfer scenario. Zero-shot corresponds to source-only calibration, whereas site-specific uses 10% target-domain calibration. Coverage and singleton values are reported in %. 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	92.3\pm2.0	1.08 \pm 0.05	92.3 \pm 4.9	96.8 \pm 0.6	1.29\pm0.11	75.7\pm9.4
	Dynamic CNN	92.0 \pm 1.9	1.03 \pm 0.03	96.4 \pm 2.5	97.4\pm0.7	1.36 \pm 0.06	67.6 \pm 6.5
	NA-DyCNN	91.6 \pm 2.1	1.02\pm0.03	97.3\pm2.1	96.9 \pm 0.4	1.31 \pm 0.03	72.4 \pm 2.5
	MobileNetV3	85.6 \pm 2.3	1.19 \pm 0.03	82.2 \pm 3.2	96.5 \pm 1.7	1.55 \pm 0.07	55.1 \pm 10.6
	ResNet18 (ref)	87.8 \pm 6.2	1.02 \pm 0.02	96.9 \pm 1.8	95.2 \pm 0.8	1.24 \pm 0.04	78.4 \pm 4.2
DRUSEN	Basic CNN	60.2 \pm 28.1	1.38 \pm 0.21	62.4 \pm 20.0	95.1 \pm 0.9	1.95 \pm 0.09	17.4 \pm 4.4
	Dynamic CNN	19.9 \pm 10.2	1.06\pm0.04	90.2\pm3.2	95.6 \pm 1.5	1.99 \pm 0.12	12.1 \pm 5.1
	NA-DyCNN	39.3 \pm 20.6	1.11 \pm 0.18	84.2 \pm 14.5	94.9 \pm 1.2	1.88\pm0.12	19.8\pm6.2
	MobileNetV3	82.3\pm7.0	1.54 \pm 0.09	48.1 \pm 10.0	95.7\pm0.5	2.14 \pm 0.27	16.3 \pm 4.5
	ResNet18 (ref)	41.0 \pm 14.7	1.08 \pm 0.05	90.0 \pm 5.8	94.1 \pm 1.4	1.75 \pm 0.01	28.8 \pm 3.8
NORMAL	Basic CNN	95.1 \pm 1.2	1.36 \pm 0.29	64.6 \pm 28.6	95.5\pm1.2	2.07 \pm 0.05	9.2 \pm 3.0
	Dynamic CNN	96.3\pm2.9	1.03\pm0.01	96.5\pm1.1	94.5 \pm 0.7	2.11 \pm 0.18	5.1 \pm 6.1
	NA-DyCNN	94.5 \pm 5.0	1.10 \pm 0.16	90.1 \pm 15.3	94.5 \pm 0.7	1.94\pm0.25	17.4\pm13.7
	MobileNetV3	92.0 \pm 5.9	1.67 \pm 0.14	36.3 \pm 16.5	94.8 \pm 1.4	2.51 \pm 0.43	3.0 \pm 2.1
	ResNet18 (ref)	95.8 \pm 5.2	1.02 \pm 0.02	97.5 \pm 2.0	94.8 \pm 0.4	1.95 \pm 0.05	9.1 \pm 6.3

Table S2. Per-class Mondrian conformal prediction performance for the $D_2 \rightarrow D_1$ transfer scenario. Zero-shot corresponds to source-only calibration, whereas site-specific uses 10% target-domain calibration. Coverage and singleton values are reported in %. 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	96.9 \pm 0.6	1.35 \pm 0.09	68.3 \pm 7.1	95.1 \pm 0.3	1.27 \pm 0.08	73.9 \pm 6.6
	Dynamic CNN	96.5 \pm 0.3	1.20 \pm 0.02	80.1 \pm 1.8	95.1 \pm 0.1	1.16 \pm 0.07	84.2 \pm 6.3
	NA-DyCNN	97.1\pm0.4	1.12\pm0.02	87.8\pm1.5	95.4\pm0.2	1.10\pm0.03	90.4\pm3.0
	MobileNetV3	96.6 \pm 0.5	1.26 \pm 0.07	75.7 \pm 6.4	95.3 \pm 0.4	1.17 \pm 0.04	83.6 \pm 3.5
	ResNet18 (ref)	96.2 \pm 0.9	1.04 \pm 0.02	95.6 \pm 1.5	94.8 \pm 0.3	1.01 \pm 0.03	97.8 \pm 2.2
DME	Basic CNN	95.2 \pm 0.2	1.29 \pm 0.07	75.2 \pm 4.9	94.7 \pm 0.7	1.19 \pm 0.05	82.4 \pm 3.4
	Dynamic CNN	94.7 \pm 1.4	1.14 \pm 0.12	86.4 \pm 11.4	94.4 \pm 0.1	1.11 \pm 0.12	89.2 \pm 10.6
	NA-DyCNN	94.8 \pm 1.0	1.07\pm0.02	93.4\pm1.5	95.0\pm0.2	1.04\pm0.02	96.0\pm1.4
	MobileNetV3	96.7\pm0.7	1.15 \pm 0.05	87.3 \pm 4.1	94.6 \pm 0.5	1.09 \pm 0.03	92.2 \pm 2.5
	ResNet18 (ref)	95.3 \pm 0.2	1.02 \pm 0.03	96.7 \pm 1.8	94.8 \pm 0.5	1.00 \pm 0.03	97.5 \pm 0.8
DRUSEN	Basic CNN	96.5 \pm 0.6	1.45 \pm 0.20	59.6 \pm 15.6	94.8 \pm 0.8	1.32 \pm 0.15	69.8 \pm 12.4
	Dynamic CNN	96.0 \pm 1.3	1.20 \pm 0.09	80.7 \pm 8.5	95.0 \pm 0.4	1.14 \pm 0.10	85.9 \pm 9.1
	NA-DyCNN	96.1\pm0.2	1.16\pm0.05	84.4\pm4.6	95.3\pm1.4	1.09\pm0.02	90.7\pm2.3
	MobileNetV3	96.6 \pm 0.3	1.24 \pm 0.05	77.2 \pm 4.5	94.5 \pm 1.5	1.16 \pm 0.04	84.5 \pm 4.1
	ResNet18 (ref)	96.2 \pm 0.2	1.04 \pm 0.01	96.0 \pm 0.8	94.5 \pm 1.2	1.00 \pm 0.02	98.2 \pm 1.2
NORMAL	Basic CNN	96.9\pm0.1	1.43 \pm 0.19	62.9 \pm 14.3	94.7 \pm 0.4	1.29 \pm 0.15	73.2 \pm 13.1
	Dynamic CNN	96.3 \pm 0.6	1.13 \pm 0.13	87.7 \pm 11.4	95.2 \pm 0.8	1.11 \pm 0.17	88.7 \pm 14.4
	NA-DyCNN	96.5 \pm 0.3	1.05\pm0.02	94.9\pm1.4	95.4\pm0.3	1.03\pm0.02	96.7\pm1.7
	MobileNetV3	95.3 \pm 0.3	1.31 \pm 0.12	72.0 \pm 10.3	95.0 \pm 0.2	1.14 \pm 0.06	86.5 \pm 5.3
	ResNet18 (ref)	94.4 \pm 1.1	1.01 \pm 0.04	97.2 \pm 2.3	95.3 \pm 0.2	0.99 \pm 0.02	97.7 \pm 0.4

Table S3. Per-class Mondrian conformal prediction performance for the $D_2 \rightarrow D_3$ transfer scenario. Zero-shot corresponds to source-only calibration, whereas site-specific uses 10% target-domain calibration. Coverage and singleton values are reported in %. 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	95.7 \pm 1.6	1.30 \pm 0.01	73.9 \pm 1.6	96.9 \pm 1.1	1.64 \pm 0.15	52.9\pm8.3
	Dynamic CNN	97.9\pm2.6	1.22 \pm 0.16	79.1 \pm 14.6	97.0\pm1.7	1.57 \pm 0.18	51.8 \pm 8.9
	NA-DyCNN	96.5 \pm 3.2	1.11\pm0.01	89.5\pm1.1	96.5 \pm 1.3	1.55\pm0.12	51.2 \pm 12.6
	MobileNetV3	89.9 \pm 9.7	1.38 \pm 0.19	70.2 \pm 14.9	95.8 \pm 1.6	1.90 \pm 0.15	35.7 \pm 14.0
	ResNet18 (ref)	93.1 \pm 2.7	1.05 \pm 0.04	94.7 \pm 4.1	96.4 \pm 0.8	1.24 \pm 0.06	80.8 \pm 4.6
DRUSEN	Basic CNN	69.3\pm22.1	1.79 \pm 0.05	33.5 \pm 5.0	95.5 \pm 0.7	2.21 \pm 0.15	12.4 \pm 3.6
	Dynamic CNN	66.2 \pm 22.1	1.74 \pm 0.04	34.6 \pm 4.0	95.1 \pm 2.1	1.93 \pm 0.32	26.8\pm15.6
	NA-DyCNN	56.6 \pm 12.2	1.49\pm0.06	53.6\pm5.5	96.3\pm1.3	1.92\pm0.26	25.6 \pm 16.7
	MobileNetV3	63.5 \pm 21.7	1.77 \pm 0.10	37.4 \pm 8.3	95.0 \pm 0.3	2.45 \pm 0.07	4.4 \pm 0.7
	ResNet18 (ref)	68.9 \pm 9.6	1.21 \pm 0.14	79.4 \pm 13.5	95.4 \pm 1.1	1.71 \pm 0.19	37.7 \pm 14.8
NORMAL	Basic CNN	87.6 \pm 8.4	1.83 \pm 0.21	33.5 \pm 12.7	95.0 \pm 1.2	2.44 \pm 0.17	9.4 \pm 4.8
	Dynamic CNN	87.4 \pm 15.1	1.70 \pm 0.24	40.1 \pm 14.8	95.6\pm0.3	2.13 \pm 0.44	20.6\pm16.6
	NA-DyCNN	93.6\pm3.9	1.42\pm0.15	59.9\pm15.2	95.1 \pm 0.2	2.20\pm0.20	12.5 \pm 8.5
	MobileNetV3	91.0 \pm 5.7	1.73 \pm 0.16	42.0 \pm 11.7	95.5 \pm 0.9	2.53 \pm 0.11	5.0 \pm 1.6
	ResNet18 (ref)	87.5 \pm 5.8	1.22 \pm 0.22	78.3 \pm 21.9	94.7 \pm 0.7	1.96 \pm 0.09	18.3 \pm 4.8

Across all transfer settings, the class-wise results are consistent with the aggregate findings reported in the main manuscript. For the relatively mild shifts (e.g., $D_1 \rightarrow D_2$ and $D_2 \rightarrow D_1$), site-specific calibration maintains coverage close to the nominal target while preserving compact prediction sets and high singleton rates. In contrast, the more challenging transfers involving D_3 exhibit substantial zero-shot under-coverage. After calibration using only 10% labeled target data,

coverage is restored to approximately 95%, but this is often accompanied by larger prediction sets and reduced singleton rates. This trade-off reflects the increased uncertainty required to recover valid conformal guarantees under severe domain shift. Across all settings, NA-DyCNN generally maintains smaller prediction sets and higher singleton rates than the other lightweight baselines while achieving comparable coverage, indicating more efficient uncertainty estimates under cross-domain deployment.