

Appendix

We provide additional insights into our study through various sections in the appendix, which include results for both versions of the VMamba models [29]: $v2$ (reported in the main paper) and $v0$. Results that do not specify the VMamba version should be assumed to refer to $v2$ by default.

In Section A, we provide detailed results on *Information Drop* experiments. Section A.1 explores information drop along the scanning direction. Sections A.2 and A.3 offer additional detailed results on Salient and Non-Salient Patch Drop and Random Patch Drop, respectively. Section A.4 presents further results on patch shuffling. Section B provides results on *Image Corruptions*. In Section B.1, robustness against global corruptions such as common corruption, and out-of-distribution datasets is evaluated. Additionally, in Section B.2, we provide model calibration results on in-distribution and out-of-distribution datasets. In Section B.3, we report results on fine-grained corruptions across different models. In Section B.4, we report image corruption results for the task of object detection and semantic segmentation. Finally, we expand our analysis of *Adversarial Attacks* across models in Section C. We present results for white-box attacks, transfer-based black-box attacks, and frequency-based attacks.

A. Robustness against Information Drop

A.1. Information Drop along the Scanning Axis

We investigate the models' behavior when information is dropped along the scanning directions. We consider three settings: (1) linearly increasing the amount of information dropped in each patch along the scanning direction, with the most information dropped in patches that are traversed last by the scanning operation, (2) dropping most of the information in the center of the scanning directions while preserving most of the information in patches that are traversed at the beginning and the end, and (3) sequentially dropping patches along the scanning directions. In Figure 5 we show qualitative samples for showing information drop along scanning directions. Figure 6 we report results across the first experiment setting across all the scanning directions. In Figure 7 we report results for the second experiment setting. Furthermore, in Figure 8 we report results across the third experiment setting across all the scanning directions.

A.2. Random Patch Drop

In Table 6 we expand the random patch drop experiments from the main paper (Table 1) to patch sizes 56×56 and 224×224 . Furthermore, in Figure 9, 10, 11, 12, 13, and 14 we expand our analysis on random patch drop to several

other models.

A.3. Salient and Non-Salient Patch Drop

In Figure 15 and 16, we report results on salient and non-salient patch drop of information.

A.4. Patch Shuffling

In Figure 17, we expand the patch shuffling experiment from the main paper to several CNN and transformer-based architectures.

B. Robustness against Image Corruptions

B.1. Robustness against Global Corruptions

In Table 7 relative corruption error is reported across models discussed in the main paper. In Figure 18 and 19, we present the relative corruption error and mean corruption error (mCE) of various models subjected to all 19 corruption methods applied to the ImageNet dataset. The common corruption consists of various types of synthetic corruptions to assess the models' robustness against image distortions like noise, blur, and compression artifacts. In Table 8 we report results on *ImageNetV2* [42], *ImageNet-A* [20], *ImageNet-R* [19], and *ImageNet-S* [11]. These datasets are designed to evaluate the robustness and generalization capabilities of models trained on the original ImageNet dataset. ImageNetV2 [42] tests the models' performance on previously unseen images that follow the same distribution as the training data, while ImageNet-A [20] contains naturally occurring adversarial examples that are difficult for models to classify correctly. ImageNet-R [19] consists of images with different artistic renditions to test the models' ability to generalize to different visual domains and styles, and ImageNet-S [11] consists of sketch-based images.

B.2. Model Calibration

Model calibration assesses how well a model's predicted confidence aligns with its actual accuracy. For example, if a model predicts a confidence level of 70% for its predictions, a well-calibrated model should have an actual accuracy close to 70%. To quantify this alignment, we use the Expected Calibration Error (ECE). ECE involves dividing predictions into M bins based on their confidence levels (e.g., 60%-70%, 70%-80%). For each bin, the average accuracy and confidence are computed, and the ECE is the weighted average of the differences between these values. Calibration can also be evaluated visually using reliability diagrams, which plot predicted confidence against actual accuracy; a well-calibrated model should show points near the diagonal. Additionally, confidence histograms reveal the distribution of prediction confidences. Evaluation is performed on both in-distribution data (e.g., ImageNet, ImageNetV2) and out-of-distribution data (e.g., ImageNet-R, ImageNet-S, ImageNet-A), with $M = 15$ bins used in all experiments.

In Figure 20, 21, and 22, we report the ECE plot the reliability diagrams for Tiny, Small, and Base version of different architectures, respectively. Figure 23, we report the ECE score. We observe that while ViT models report a low ECE error on in-distribution datasets, the error increases significantly for out-of-distribution datasets. An exception is the Base model, which can be attributed to the ViT-B model being pretrained on a larger dataset (ImageNet-21k) compared to other models. For the hybrid-based VSSM model MambaVision we observe a contrasting behaviour with ViT models, while the ECE score on in-distribution datasets is relatively high, it improves significantly on out-of-distribution datasets. The improvement in ECE for MambaVision on out-of-distribution datasets suggests that the hybrid approach of combining vision transformers with other model components may enhance calibration performance in these scenarios.

B.3. Robustness against Fine-grained Corruptions

In Table 9 (*left*) we report results on ImageNet-E dataset across several CNNs, Transformers, and VSSM-based models, and in Table 9 (*right*), we report results on ImageNet-B dataset.

B.4. Robustness Evaluation for Object Detection and Semantic Segmentation

In Table 10, we report AP scores of different architectures on COCO-DC dataset. On COCO-DC, we observe that VMamba models’ high performance on clean images does not translate to color and texture background variations. Swin-S model achieves the highest average AP score of 56.70 across the background variations, followed by score of 56.20 by VMamba-S model. Figure 24 shows the mIoU, mAcc, and aAcc scores on ADE-20K after applying common corruptions to the dataset. Similarly, Figure 25 displays the mAP, APs, APm, and API scores on the COCO-C dataset. All the scores are averaged across all severity levels of each corruption.

C. Robustness against Adversarial Attacks

In this section, we expand our analysis on adversarial attacks and their transferability across different models. We include more model families, such as DeiT, DenseNet, and VGG. Figures 26, 27, and 28 present the robust accuracy of various models under both white-box and black-box settings for FGSM, PGD, and MIFGSM attacks, respectively. All adversarial examples are crafted with a perturbation budget of $\epsilon = \frac{8}{255}$. For PGD and MIFGSM attacks, we use 20 iterations to craft the adversarial examples.

In Table 11, we evaluate the robustness of VMamba and MambaVision against frequency-specific adversarial attacks crafted using Projected Gradient Descent (PGD). These perturbations are constrained to designated frequency

bands through a discrete cosine transform (DCT) mask filter [34]. Results in Tab. 11 (*left*) demonstrate that VMamba and MambaVision maintain robustness above 90% for low-frequency perturbations up to $\epsilon = 16$, indicating strong resilience similar to the ConvNext and Swin transformer counterparts. ViT models exhibit the most performance drop under low-frequency adversarial attacks. For high-frequency adversarial attacks (Tab. 11 (*middle*)), the robustness of all models decreases more rapidly with increasing perturbation strength, although ViT-based models show the highest robustness. Finally, for standard attacks where the complete frequency range is used to generate adversarial examples, VSSM models display higher robustness compared to other models, including ConvNext, ViT, and Swin.

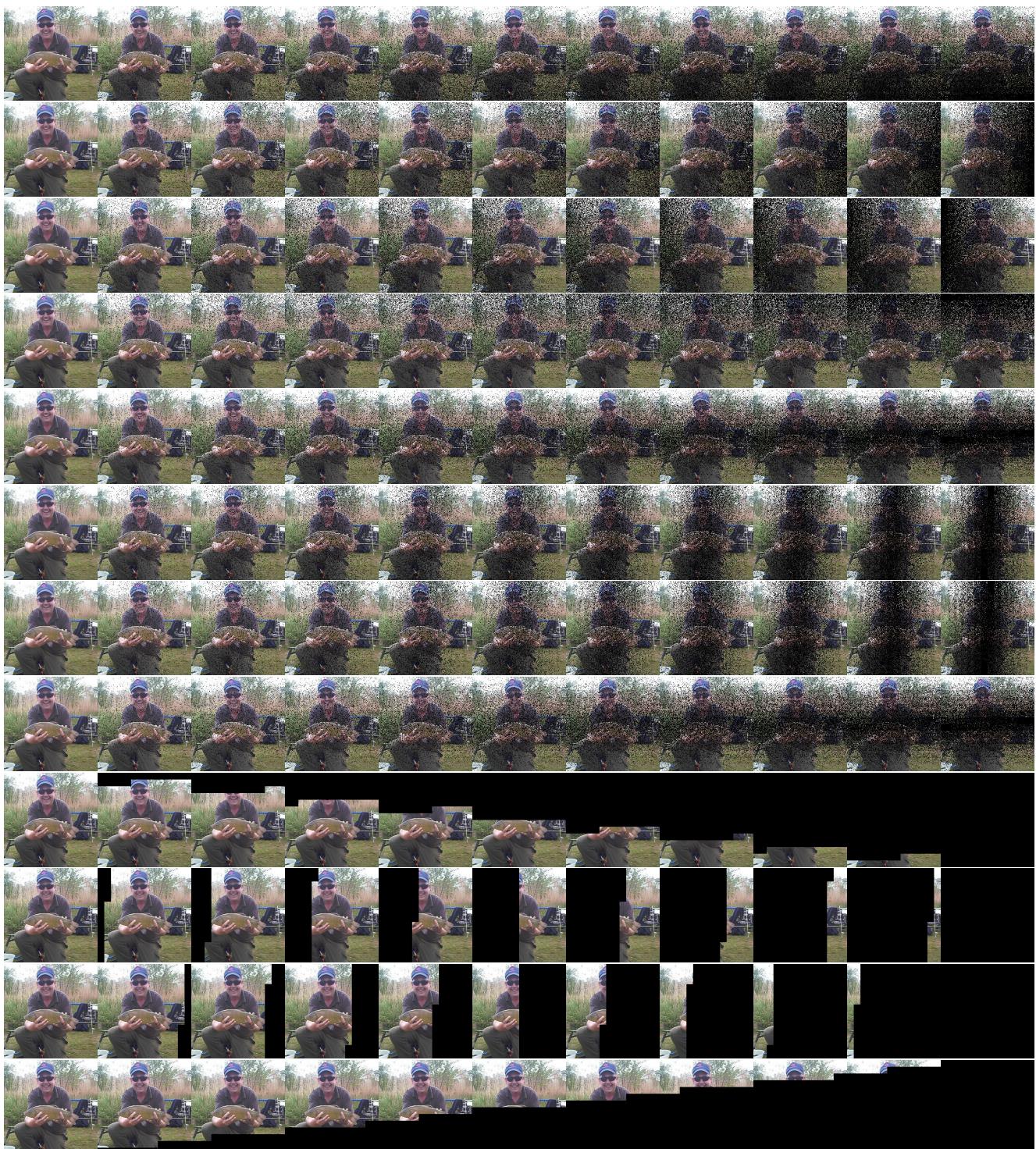


Figure 5. Information drop along scanning direction: The top four rows represent linearly increasing information drop along the four scanning directions at patch size 14×14 . The center four rows represent linearly increasing information drop till the center along the four scanning directions at patch size 14×14 . The bottom four rows represent sequentially dropping patches along the four scanning directions at patch size 14×14 .

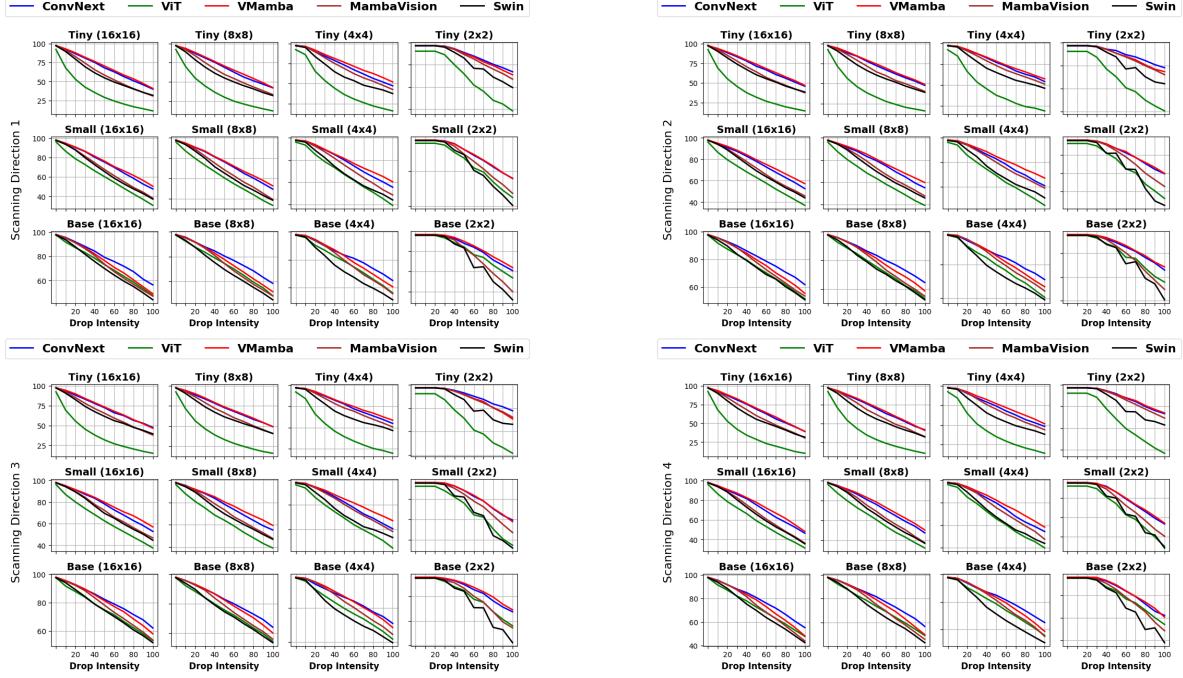


Figure 6. Information drop of Tiny and Small family of models along the scanning direction: the image is split into a sequence of fixed-size non-overlapping patches of size 16x16, 8x8, 4x4, and 2x2. We report results of linearly increasing the number of pixels dropped from each patch to the maximum threshold (*Drop Intensity*) along the scanning direction. Top row shows results for top-to-bottom (*Direction 1*) and left-to-right (*Direction 2*). Bottom row shows results for right-to-left (*Direction 3*) and bottom-to-top (*Direction 4*).

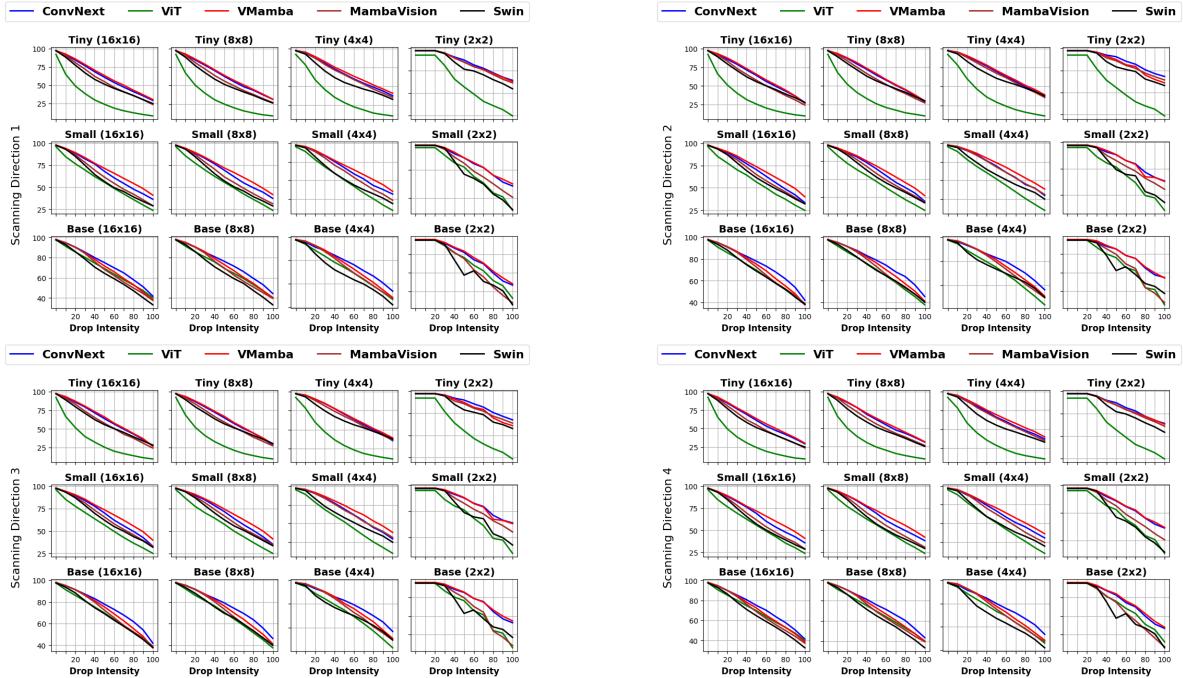


Figure 7. Information drop of Tiny and Small family of models along the scanning direction: the image is split into a sequence of fixed-size non-overlapping patches of size 16x16, 8x8, 4x4, and 2x2. We report results of linearly increasing the number of pixels dropped from each patch to the maximum threshold (*Drop Intensity*) at the center of the scanning direction and then again linearly decreased till the end. The top row shows results for top-to-bottom (*Direction 1*) and left-to-right direction (*Direction 2*). The bottom row shows results for right-to-left (*Direction 3*) and bottom-to-top direction (*Direction 4*).

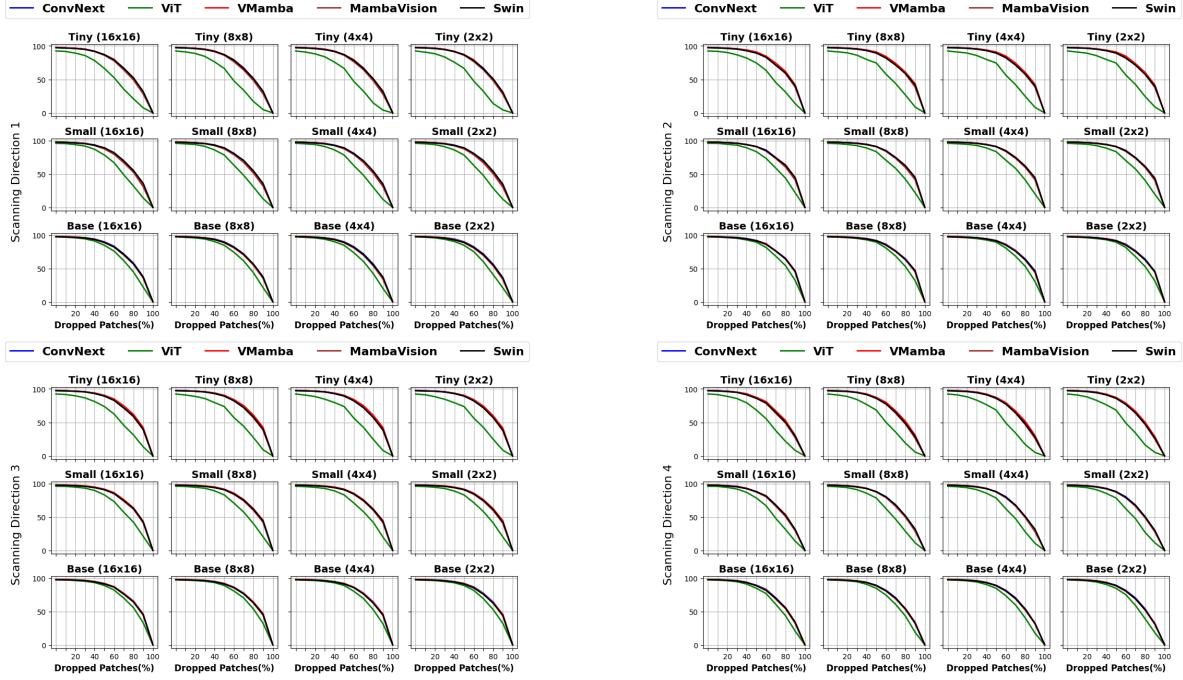


Figure 8. Information drop of Tiny and Small family of models along the scanning direction: the image is split into a sequence of fixed-size non-overlapping patches of size 16x16, 8x8, 4x4, and 2x2. We report results of sequentially dropping patches along the scanning direction. The top row shows results for top-to-bottom (*Direction 1*) and left-to-right (*Direction 2*). The bottom row shows results for right-to-left (*Direction 3*) and bottom-to-top (*Direction 4*).

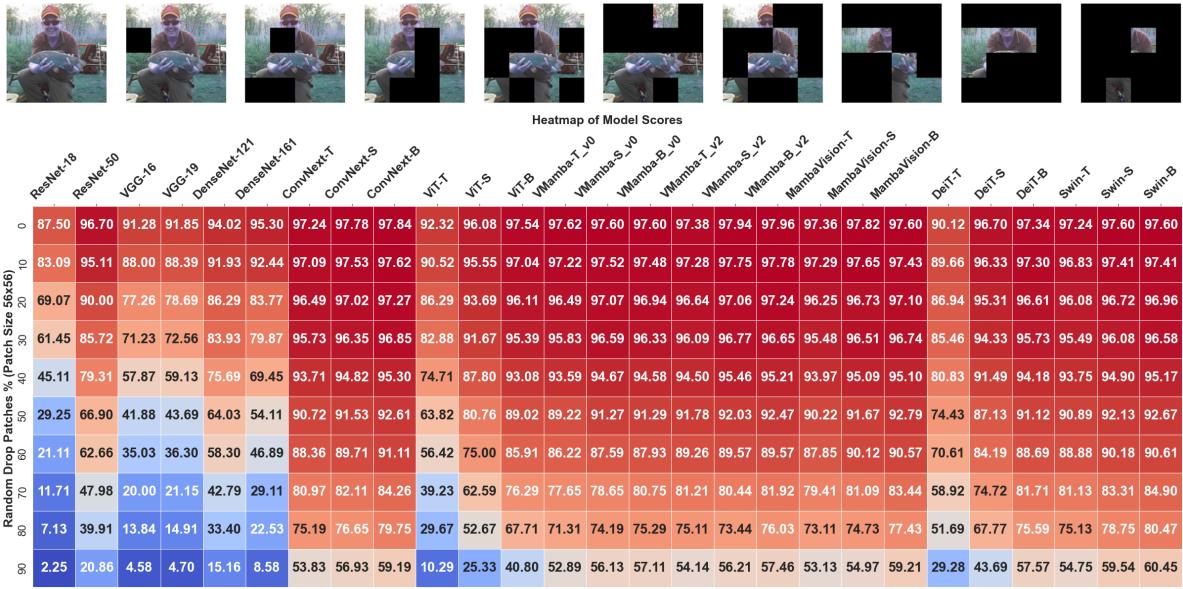


Figure 9. Top-1 classification accuracy of various architectures under random patch drop occlusions, using 56 × 56 patch size.

Table 6. Top-1 classification accuracy reported random patch drop occlusion using 16×16 , 8×8 , 4×4 and 1×1 patch sizes.

ResNet-50	ConvNext-T	ConvNext-S	ConvNext-B	ViT-T	ViT-S	ViT-B	VMamba-T	VMamba-S	VMamba-B	MambaVision-T	MambaVision-S	MambaVision-B	Swin-T	Swin-S	Swin-B
Patch Size 16×16 (Percentage of patch drop increasing from top to bottom (10% to 90%))															
96.70	97.24	97.78	97.84	92.30	96.08	97.54	97.38	97.94	97.96	97.36	97.82	97.60	97.24	97.60	97.60
75.27	96.49	97.19	97.37	90.83	95.39	96.85	96.49	96.61	97.25	96.24	96.80	97.09	96.76	97.38	97.32
39.93	94.63	95.27	96.48	88.09	94.29	96.27	95.16	92.97	96.25	92.91	94.67	95.74	96.11	96.84	96.79
17.91	89.99	91.12	95.29	85.26	92.35	95.08	93.45	89.74	95.21	87.14	91.19	93.25	94.88	95.88	96.17
6.73	81.43	84.63	93.03	80.08	90.15	92.78	90.52	84.82	93.46	77.02	84.56	88.05	93.38	94.35	95.03
2.43	70.07	74.44	88.76	72.49	85.10	89.21	86.52	78.41	90.89	61.79	76.07	78.92	91.05	92.21	93.25
1.05	57.59	60.15	82.35	61.34	76.56	82.08	80.52	68.40	87.03	42.16	62.44	63.89	87.96	87.84	90.43
0.56	44.67	44.71	71.25	45.63	62.63	70.25	70.39	52.72	79.96	20.73	42.16	43.33	80.65	79.71	84.70
0.45	31.29	28.82	57.71	25.86	41.73	50.08	56.23	34.51	67.56	5.62	21.07	22.21	70.37	66.60	74.38
0.43	16.73	14.98	33.67	7.85	15.86	19.68	34.83	16.82	41.55	1.95	7.08	11.08	47.16	47.85	53.54
Patch Size 8×8 (Percentage of patch drop increasing from top to bottom (10% to 90%))															
96.70	97.24	97.78	97.84	92.32	96.08	97.54	97.38	97.94	97.96	97.36	97.82	97.60	97.24	97.60	97.60
44.91	86.69	91.39	95.63	70.18	88.90	94.23	87.86	85.93	90.20	89.57	91.25	92.32	96.44	96.77	96.76
12.44	68.38	81.13	90.83	42.19	79.72	88.37	79.91	78.23	84.43	73.34	79.03	84.32	95.04	94.87	95.84
3.79	55.12	68.35	84.08	16.91	65.39	76.63	70.40	70.95	78.47	51.09	60.21	71.70	92.87	92.08	94.49
1.35	39.05	54.58	73.51	4.59	46.17	60.12	57.34	59.09	70.04	29.64	39.81	55.01	90.13	88.21	92.40
0.56	23.89	37.94	58.05	1.25	25.91	39.97	43.09	44.08	58.29	13.73	23.86	34.67	85.40	81.87	88.36
0.33	13.33	21.97	40.37	0.45	11.03	21.56	28.25	27.85	43.00	4.90	11.86	16.70	78.76	71.63	81.85
0.21	5.95	9.85	21.51	0.21	3.78	8.85	14.69	12.45	25.75	1.51	4.59	5.90	68.40	53.81	70.03
0.24	2.08	2.31	7.85	0.14	1.02	2.49	5.11	2.45	10.07	0.50	1.06	1.16	52.41	29.58	49.35
0.25	0.46	0.56	1.33	0.16	0.39	0.59	0.75	0.43	1.65	0.22	0.31	0.21	28.46	8.32	23.67
Patch Size 4×4 (Percentage of patch drop increasing from top to bottom (10% to 90%))															
96.70	97.24	97.78	97.84	92.30	96.08	97.54	97.38	97.94	97.96	97.36	97.82	97.60	97.24	97.60	97.60
30.11	84.75	88.74	90.51	29.62	82.23	90.04	86.34	87.17	90.48	80.16	85.39	89.79	92.51	92.97	94.68
12.02	64.23	77.49	81.82	9.25	65.28	77.00	72.09	75.43	77.77	57.80	57.14	71.68	84.20	86.17	90.71
6.49	33.19	47.91	70.73	3.63	43.92	56.17	51.01	55.23	53.83	29.05	24.13	44.27	72.07	77.69	85.35
2.91	13.34	19.71	54.95	1.58	21.73	32.19	27.17	33.21	29.41	11.75	8.59	23.15	57.71	65.13	78.48
1.46	5.17	6.67	34.97	0.70	8.96	15.24	10.51	15.82	11.57	4.00	2.84	10.25	41.57	47.95	69.11
0.75	2.07	2.23	17.37	0.41	3.69	6.01	3.07	5.41	3.40	1.41	0.90	3.89	26.49	29.31	54.59
0.40	0.85	0.90	6.21	0.24	1.45	1.91	0.83	1.28	0.63	0.50	0.39	1.05	12.17	12.45	35.95
0.21	0.33	0.44	1.70	0.17	0.55	0.71	0.23	0.30	0.24	0.29	0.17	0.23	3.48	2.47	14.47
0.13	0.20	0.25	0.40	0.19	0.24	0.25	0.19	0.21	0.22	0.19	0.13	0.14	0.64	0.41	1.55
Patch Size 1×1 (Percentage of patch drop increasing from top to bottom (10% to 90%))															
96.70	97.24	97.78	97.84	92.32	96.08	97.54	97.38	97.94	97.96	97.36	97.82	97.60	97.24	97.60	97.60
54.40	83.43	87.12	89.68	47.56	76.30	85.64	85.87	89.55	89.73	80.38	86.96	90.11	74.40	85.10	85.14
37.39	69.17	76.25	79.47	27.12	61.71	75.93	71.95	77.99	79.11	57.76	69.25	76.01	48.34	63.30	62.51
26.21	53.81	62.42	68.38	15.67	47.69	65.85	56.63	64.54	66.81	39.91	50.91	57.74	33.44	45.94	49.32
17.80	38.14	45.85	58.50	9.02	34.97	55.16	43.11	52.29	54.38	27.55	35.95	40.89	25.06	32.73	39.90
11.85	26.13	31.79	47.28	5.11	23.24	43.27	32.61	41.20	42.08	18.56	24.83	27.92	18.84	23.16	31.09
7.19	17.81	21.99	35.81	2.78	14.30	31.58	23.08	30.97	29.57	11.42	16.06	19.34	12.75	15.41	22.05
3.94	11.18	15.01	21.95	1.39	7.82	19.59	13.39	21.11	17.70	5.97	9.18	11.85	6.69	9.58	11.99
1.79	4.47	8.49	9.15	0.61	3.51	10.46	4.65	9.64	6.79	2.19	3.19	4.87	2.15	4.15	4.05
0.62	0.99	2.91	2.62	0.29	1.23	3.13	0.79	1.47	1.23	0.71	0.63	1.08	0.61	1.08	0.95

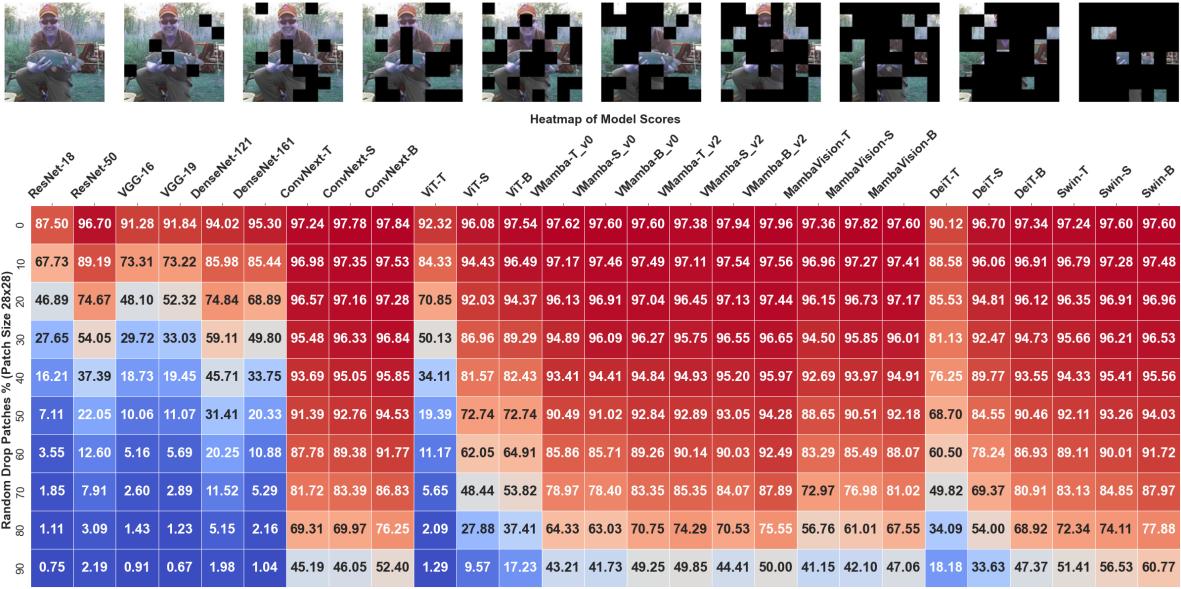


Figure 10. Top-1 classification accuracy of various architectures under random patch drop occlusions, using 28 × 28 patch size.

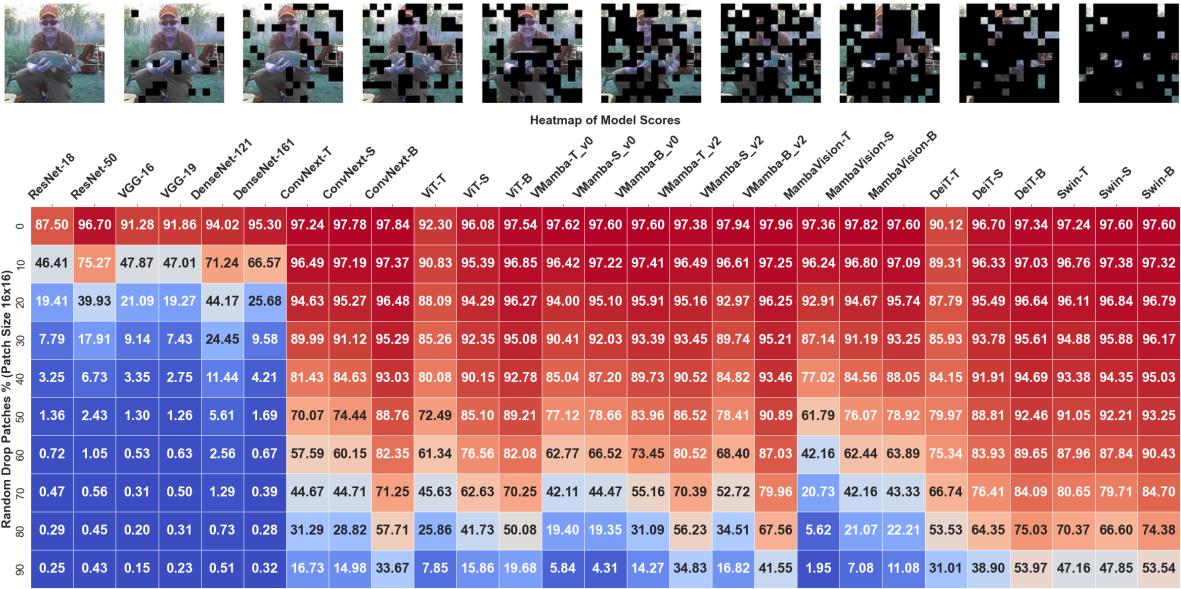


Figure 11. Top-1 classification accuracy of various architectures under random patch drop occlusions, using 16 × 16 patch size.

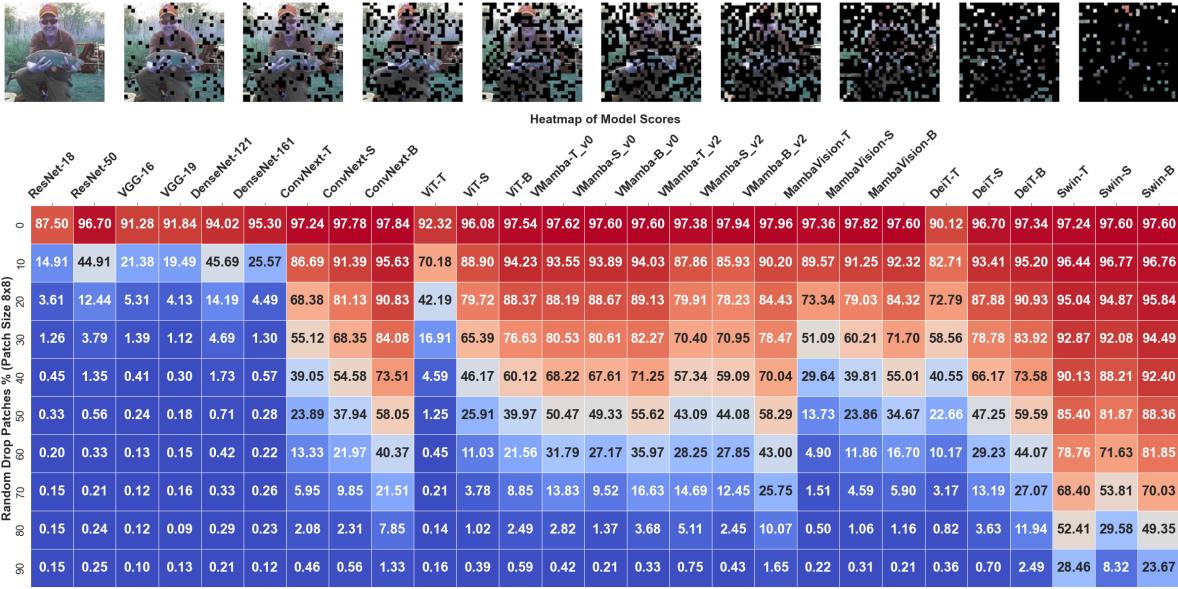


Figure 12. Top-1 classification accuracy of various architectures under random patch drop occlusions, using 8 × 8 patch size.

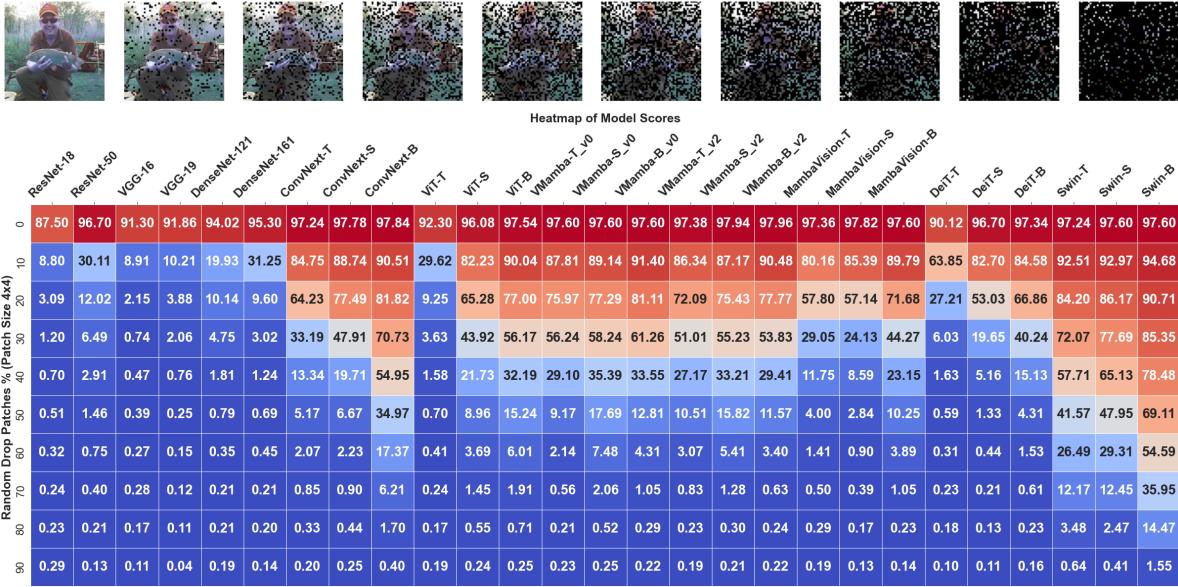


Figure 13. Top-1 classification accuracy of various architectures under random patch drop occlusions, using 4 × 4 patch size.

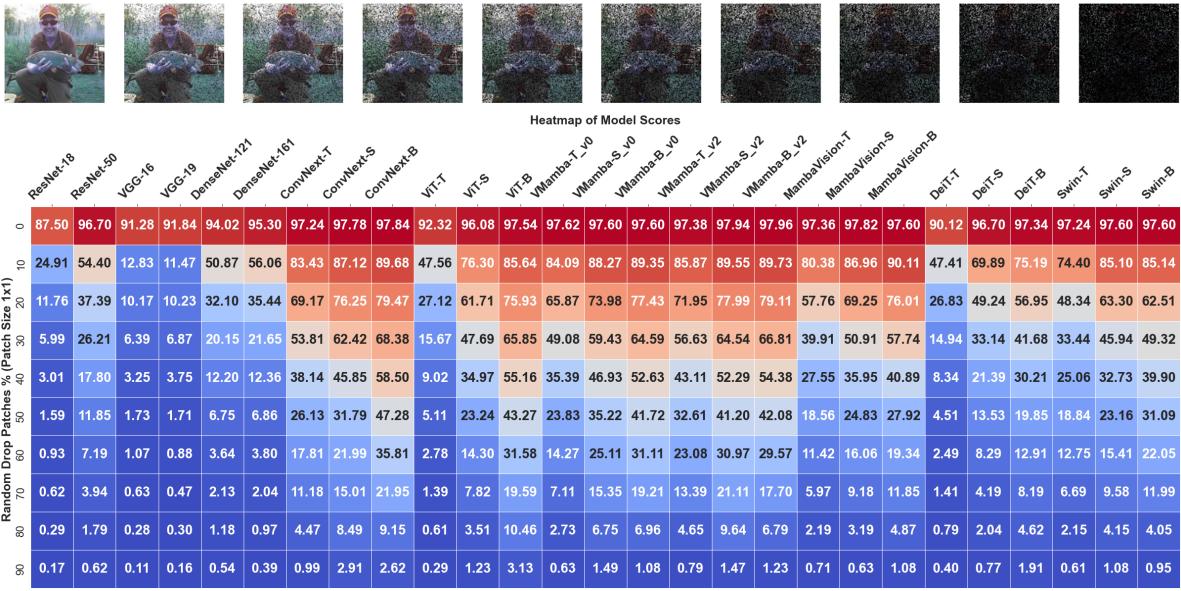


Figure 14. Top-1 classification accuracy of various architectures under random patch drop occlusions, using 1×1 patch size.

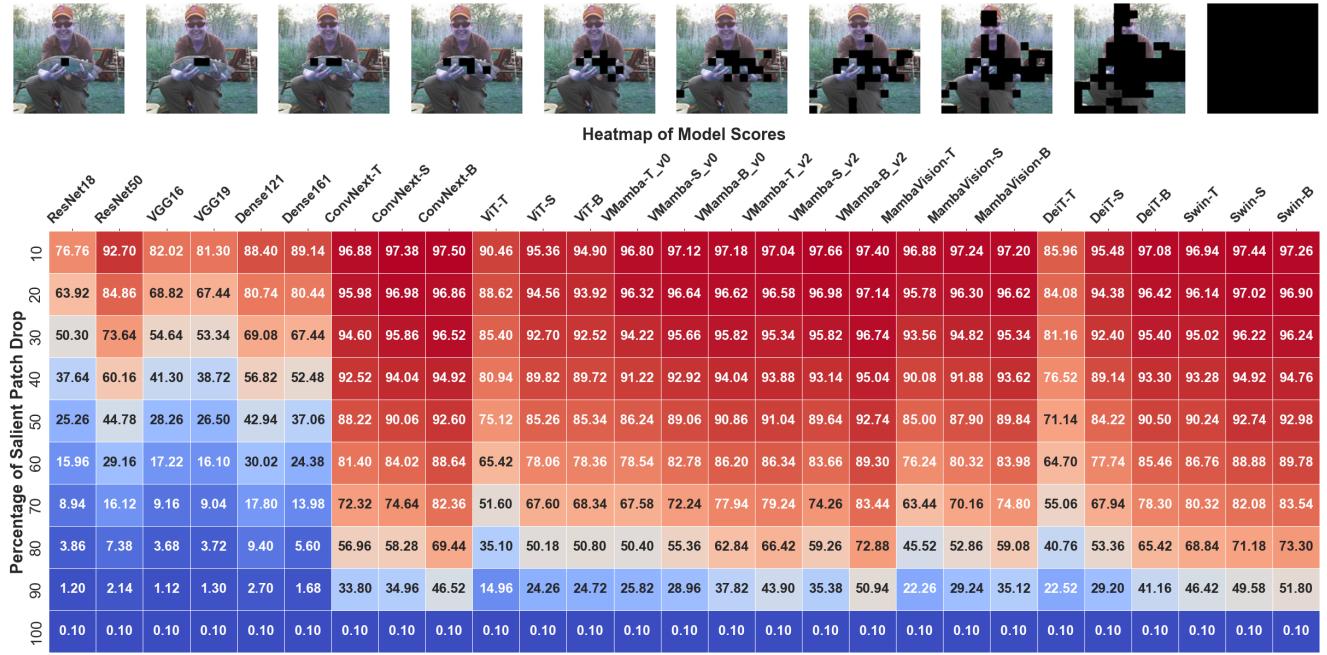


Figure 15. Top-1 classification accuracy reported under salient patch drop occlusion using 16×16 patch size.

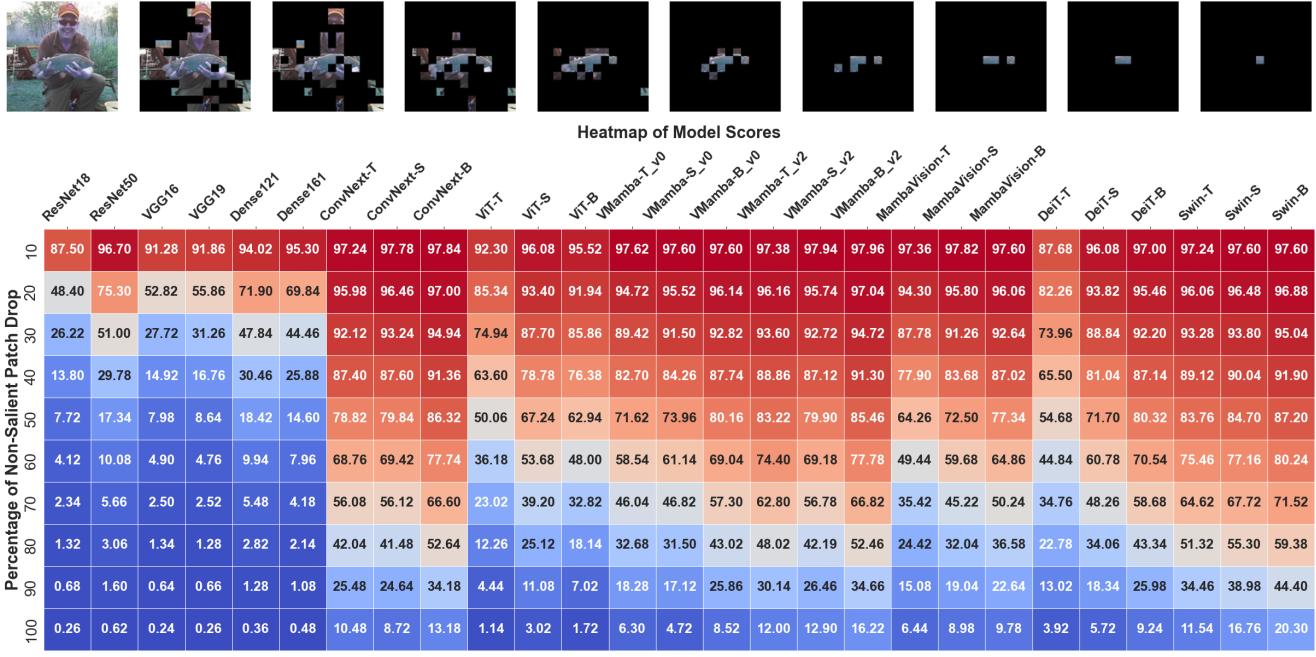


Figure 16. Top-1 classification accuracy reported under non-salient patch drop occlusion using 16×16 patch size.

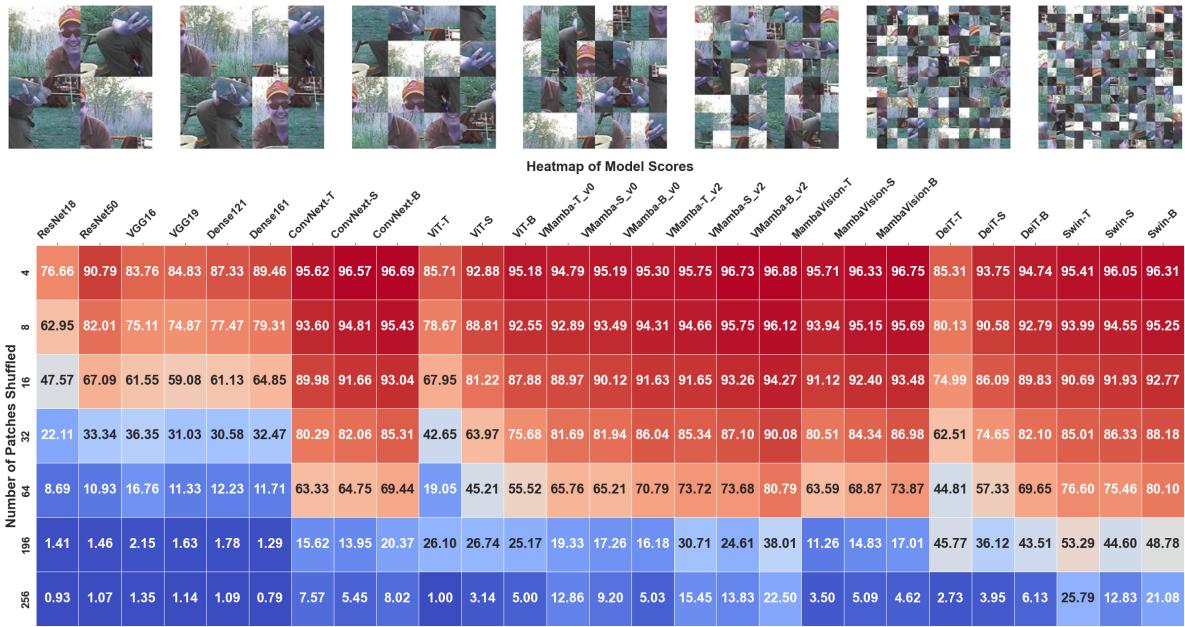


Figure 17. Top-1 classification accuracy among various architectures under increasing patch shuffling.

Table 7. Corruption Error (CE) of various architectures on ImageNet-C corruptions over multiple intensity levels. Results are relative to the CE on ResNet-50 and are evaluated on 5000 ImageNet validation set images.

Corruptions ↓	ResNet-50	ConvNext-T	ConvNext-S	ConvNext-B	ViT-T	ViT-S	ViT-B	VMamba-T	VMamba-S	VMamba-B	MambaVision-T	MambaVision-S	MambaVision-B	Swin-T	Swin-S	Swin-B							
Brightness	100.00	68.20	58.63	52.15	223.14	101.09	41.67	221.68	85.83	67.47	61.60	53.06	52.94	64.39	53.73	51.30	75.89	65.35	63.54				
Contrast	100.00	40.14	30.29	28.96	426.29	186.49	103.16	197.49	69.63	45.66	38.17	30.02	27.70	39.69	30.82	28.23	52.05	34.93	29.72	44.41	42.09	41.70	
Defocus Blur	100.00	85.70	77.82	73.87	103.27	70.64	42.81	107.14	76.91	70.41	84.96	72.32	73.00	77.80	71.71	67.04	87.50	66.82	81.49				
Elastic Transform	100.00	81.11	72.18	69.44	70.34	55.69	41.43	78.12	60.58	55.18	82.78	69.89	72.21	80.83	69.72	69.65	78.96	79.82	75.06	86.36	66.53	76.78	
Fog	100.00	113.68	86.18	91.19	312.26	119.54	68.79	233.22	102.15	73.18	87.47	64.93	55.88	62.78	40.55	36.40	62.13	46.48	39.35	84.32	87.24	68.61	
Frost	100.00	64.12	57.09	50.75	146.53	75.52	27.19	100.11	50.99	36.20	42.05	36.58	30.79	39.28	32.54	57.65	49.73	40.25	47.77	40.72	38.63		
Gaussian Blur	100.00	86.69	84.58	77.75	107.84	82.50	55.22	101.86	81.03	74.64	88.91	77.05	77.21	86.91	74.64	74.54	79.64	74.80	70.72	90.59	81.36	85.37	
Gaussian Noise	100.00	52.01	51.36	56.02	130.74	66.75	30.16	112.66	61.24	43.01	56.81	38.29	52.62	74.88	42.62	40.77	115.29	110.36	96.79	73.31	59.86	54.49	
Glass Blur	100.00	89.71	85.46	81.43	98.30	84.73	66.52	98.76	87.57	83.31	93.59	86.00	87.28	88.43	82.45	81.06	90.58	90.13	89.29	96.53	86.27	90.24	
Impulse Noise	100.00	47.18	44.71	45.78	122.95	61.19	28.27	104.68	58.78	40.71	44.55	31.93	39.72	66.31	37.76	36.25	92.95	84.21	75.73	63.83	48.63	48.68	
JPEG Compression	100.00	69.79	59.62	60.98	131.33	78.90	43.49	142.59	86.18	72.45	82.29	65.82	63.30	76.58	62.86	60.45	75.33	69.63	62.72	144.27	88.13	108.22	
Motion Blur	100.00	65.05	56.10	51.42	106.50	58.86	33.26	111.65	75.38	66.51	65.84	54.55	54.13	62.49	51.18	47.52	63.42	58.58	53.44	80.88	50.12	62.42	
Pixelate	100.00	82.35	71.52	66.00	52.11	27.31	15.88	92.04	54.87	43.55	80.34	66.30	65.02	75.51	41.46	63.78	67.42	67.31	59.34	109.24	100.70	101.51	
Saturate	100.00	63.44	53.74	47.61	250.68	116.28	48.06	201.17	87.74	64.97	57.26	47.97	47.07	65.46	49.91	48.74	60.37	52.12	48.11	75.29	62.40	58.34	
Shot Noise	100.00	52.69	51.85	56.73	128.00	66.02	30.12	107.57	63.27	46.94	56.20	38.75	51.72	70.08	40.87	38.68	105.73	104.55	91.49	73.51	62.09	57.99	
Snow	100.00	54.55	48.57	43.18	134.88	61.12	26.04	105.38	62.87	47.10	54.29	46.02	42.14	48.89	43.42	38.48	54.68	48.88	42.89	57.70	39.38	47.00	
Spatter	100.00	39.86	32.09	29.41	117.51	48.16	21.03	103.55	57.64	44.78	33.10	28.59	24.33	32.35	25.82	23.68	39.28	29.69	24.92	32.71	23.79	22.92	
Speckle Noise	100.00	44.78	40.91	34.08	132.49	60.64	26.89	107.42	57.94	43.84	42.92	30.43	34.07	52.65	34.57	32.58	84.20	79.68	64.73	61.92	48.85	48.80	
Zoom Blur	100.00	82.47	68.22	64.74	125.18	84.89	52.90	128.95	99.88	85.83	82.72	70.04	70.37	89.19	74.09	73.40	76.07	71.43	65.52	93.27	77.34	77.96	
mCE	100.00	67.55	59.52	56.92	153.7	79.28	42.26	65.05	53.08	53.76	66.58	50.50	49.61	77.86	63.04	64.98							

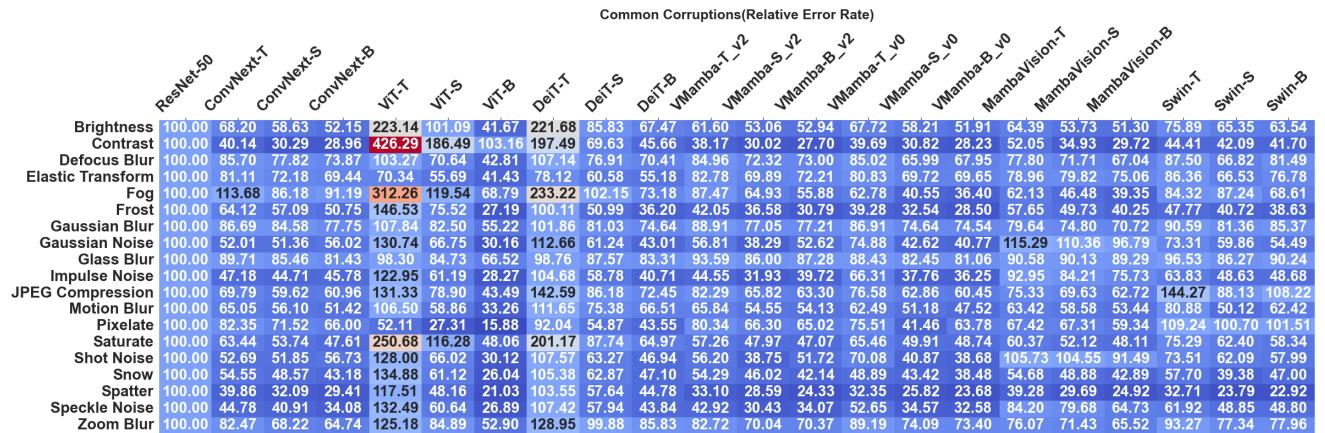


Figure 18. Corruption Error(CE) values for different corruptions and architectures on ImageNet-C, with error rates relative to ResNet-50.

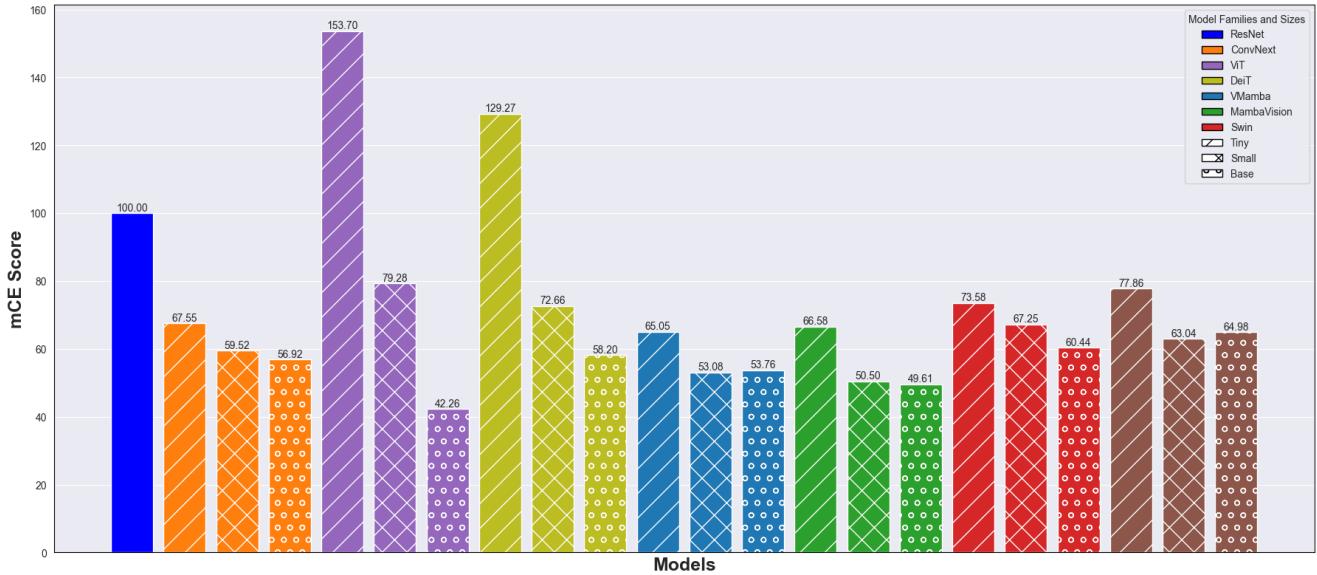


Figure 19. mCE for different corruptions and architectures on ImageNet-C, with error rates relative to ResNet-50.

Table 8. Comparison in domain generalization setting. Models trained on ImageNet are evaluated on datasets with domain shifts.

Model ↓	ImageNet	ImageNetV2	ImageNet-S	ImageNet-A	ImageNet-R	Average
ConvNext-T	81.87	70.67 _(-11.20)	33.96 _(-47.91)	10.48 _(-71.39)	32.53 _(-49.34)	36.90 _(-44.97)
ViT-T	75.35	63.05 _(-12.30)	20.88 _(-54.47)	3.31 _(-72.04)	20.29 _(-55.06)	26.88 _(-48.47)
Swin-T	80.91	69.31 _(-11.60)	29.24 _(-51.67)	8.93 _(-71.98)	28.52 _(-52.39)	34.00 _(-46.91)
VMamba-T(v0)	81.92	70.82 _(-11.10)	33.01 _(-48.91)	13.56 _(-68.36)	32.11 _(-49.81)	37.37 _(-44.55)
VMamba-T(v2)	82.28	71.16 _(-11.12)	33.99 _(-48.29)	12.08 _(-70.20)	32.05 _(-50.23)	37.32 _(-44.96)
MambaVision-T	82.10	71.50 _(-10.60)	33.49 _(-48.61)	13.39 _(-68.71)	32.28 _(-49.82)	37.66 _(-44.44)
ConvNext-S	82.82	72.07 _(-10.75)	37.16 _(-45.66)	14.43 _(-68.39)	35.57 _(-47.25)	39.81 _(-43.01)
ViT-S	81.40	69.98 _(-11.42)	32.77 _(-48.63)	13.09 _(-68.31)	31.14 _(-50.26)	36.74 _(-44.66)
Swin-S	82.90	71.62 _(-11.28)	31.97 _(-50.93)	15.72 _(-67.18)	31.93 _(-50.97)	37.81 _(-45.09)
VMamba-S(v0)	83.15	72.89 _(-10.26)	38.05 _(-45.10)	17.28 _(-65.87)	37.05 _(-46.10)	41.32 _(-41.83)
VMamba-S(v2)	83.48	73.01 _(-10.47)	36.98 _(-46.50)	16.45 _(-67.03)	35.75 _(-47.73)	40.55 _(-42.93)
MambaVision-S	83.22	72.62 _(-10.60)	35.53 _(-47.69)	15.97 _(-67.25)	33.96 _(-49.26)	39.52 _(-43.70)
ConvNext-B	83.75	73.68 _(-10.07)	38.23 _(-45.52)	18.16 _(-65.59)	36.66 _(-47.09)	41.68 _(-42.07)
ViT-B	84.40	73.84 _(-10.56)	43.01 _(-41.39)	24.09 _(-60.31)	41.03 _(-43.37)	45.49 _(-38.91)
Swin-B	83.08	72.09 _(-10.99)	32.62 _(-50.46)	17.95 _(-65.13)	33.23 _(-49.85)	38.97 _(-44.11)
VMamba-B(v0)	83.48	72.97 _(-10.51)	38.24 _(-45.24)	18.88 _(-64.60)	37.33 _(-46.15)	41.86 _(-41.62)
VMamba-B(v2)	83.76	73.22 _(-10.54)	38.53 _(-45.23)	18.35 _(-65.41)	35.99 _(-47.77)	41.52 _(-42.24)
MambaVision-B	83.96	73.84 _(-10.12)	36.69 _(-47.27)	21.69 _(-62.27)	35.72 _(-48.24)	41.98 _(-41.98)

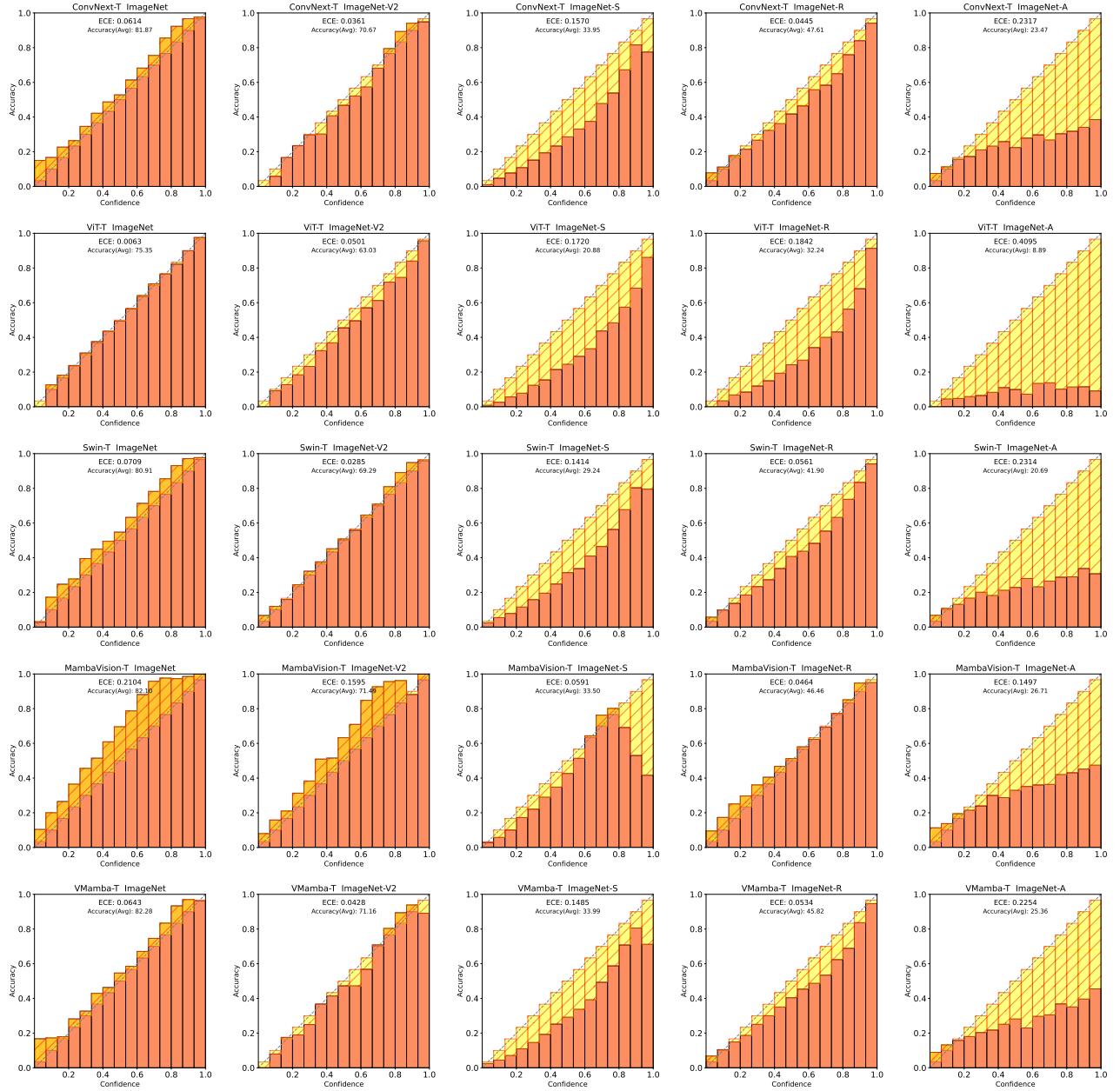


Figure 20. Calibration Results: Reliability diagrams and ECE on ConvNext-T, ViT-T, Swin-T, VMamba-T, and MambaVision-T across ImageNet, ImageNet-V2, ImageNet-S, ImageNet-R, and ImageNet-A.

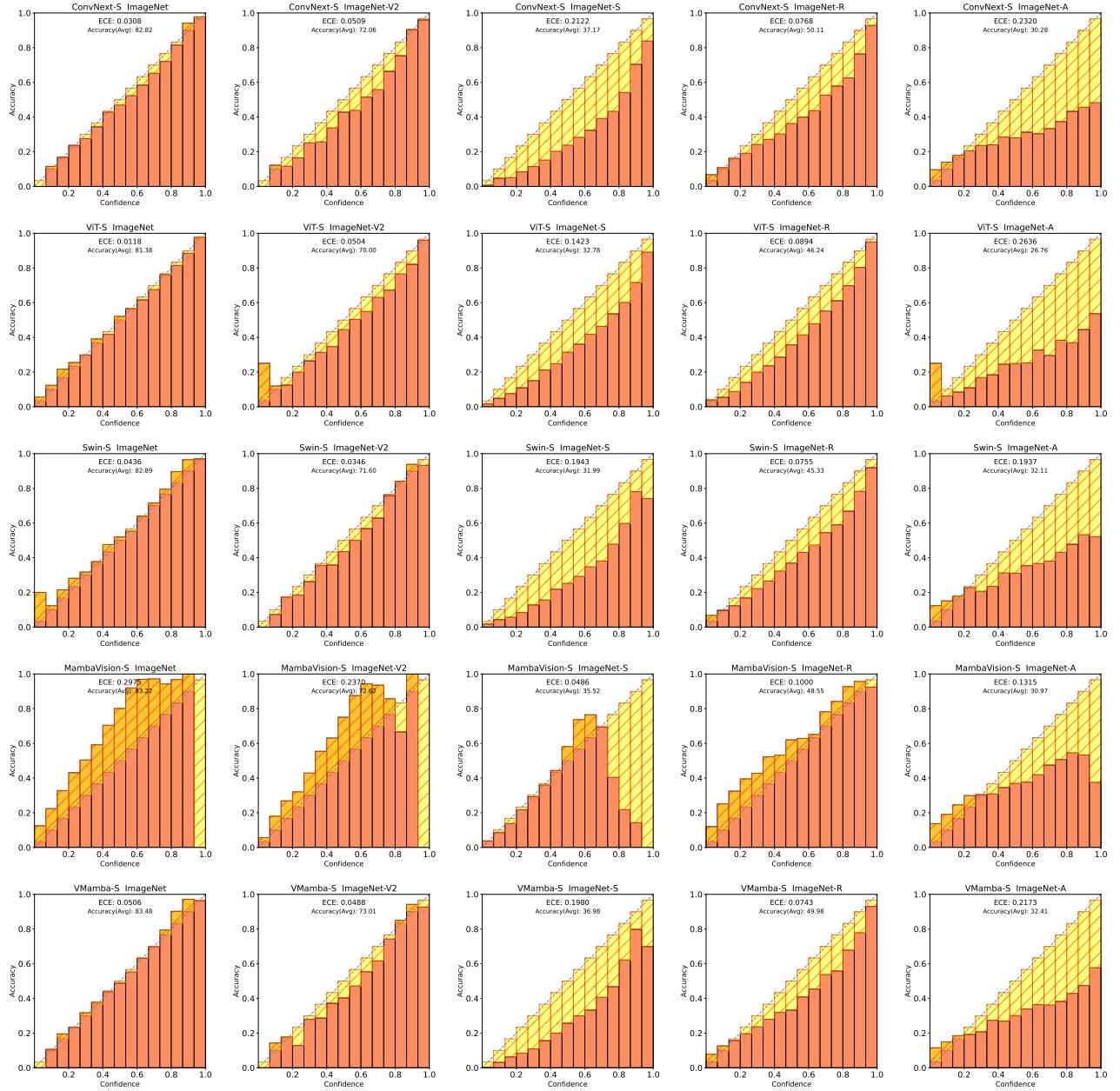


Figure 21. Calibration Results: Reliability diagrams and ECE on ConvNext-S, ViT-S, Swin-S, VMamba-S, and MambaVision-S across ImageNet, ImageNet-V2, ImageNet-S, ImageNet-R, and ImageNet-A.

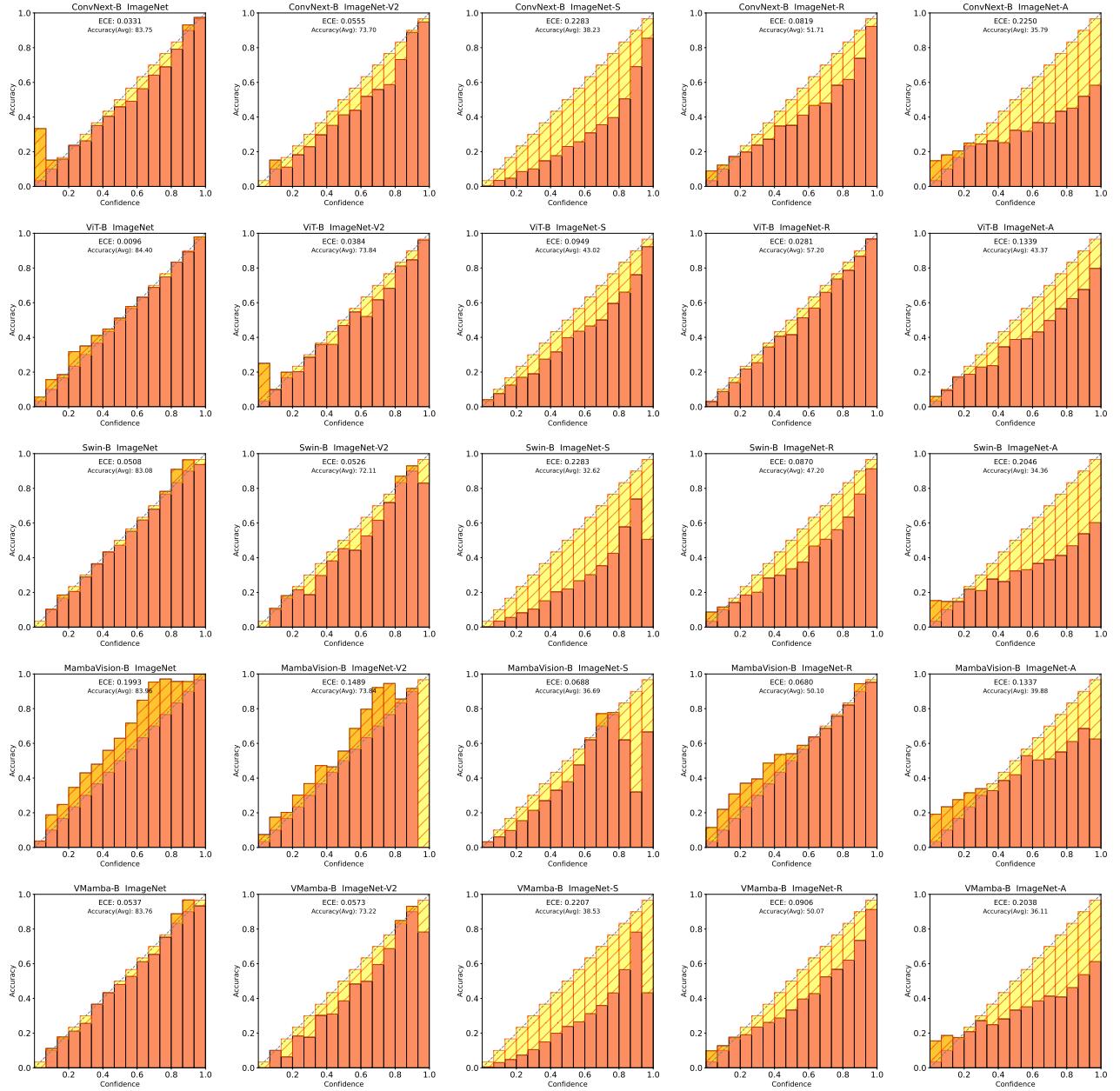


Figure 22. Calibration Results: Reliability diagrams and ECE on ViT-B, Swin-B, VMamba-B, and MambaVision-B across ImageNet, ImageNet-V2, ImageNet-S, ImageNet-R, and ImageNet-A.

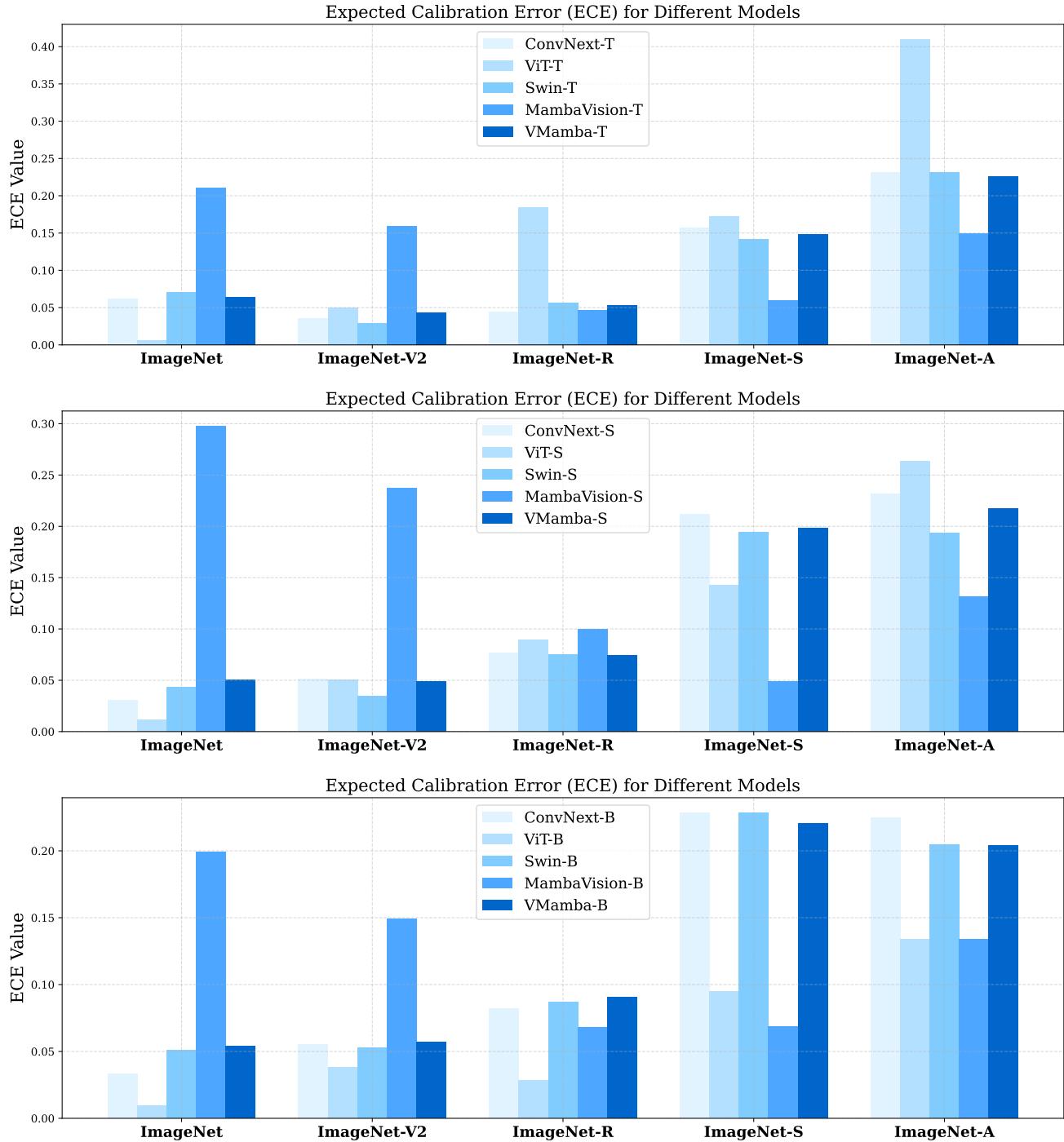


Figure 23. ECE error across classification models across ImageNet, ImageNet-V2, ImageNet-S, ImageNet-R, and ImageNet-A.

Table 9. Top-1 classification accuracy of various architectures on the ImageNet-E dataset [27] (*left*) and ImageNet-B dataset [35] (*right*).

Dataset →	ImageNet-E								ImageNet-B					
Model ↓	$\lambda = -20$	$\lambda = 20$	$\lambda = 20(\text{adv})$	Random-BG	0.1	0.08	0.05	Random Pos.	Original	Original	Caption	Class	Color	Texture
ResNet-50	88.74	86.76	73.02	84.05	89.19	86.60	77.34	73.30	94.55	98.60	94.00	96.60	88.20	85.70
VGG16	84.14	79.62	62.59	77.27	83.63	80.16	70.93	64.15	91.06	94.10	88.20	93.70	74.80	75.00
VGG19	83.89	80.15	63.21	77.80	81.16	81.01	70.54	64.91	92.05	94.50	87.80	93.30	77.10	77.80
DenseNet-121	84.93	81.78	62.36	79.55	85.15	81.17	70.08	64.50	92.78	96.20	90.20	95.10	81.10	80.00
DenseNet-161	87.48	85.48	67.72	82.79	87.22	84.41	75.00	70.08	93.04	97.50	90.70	94.70	81.10	80.10
ConvNext-T	90.95	90.03	76.88	88.09	93.01	90.87	83.09	80.19	96.09	98.20	93.20	95.10	88.80	87.40
ConvNext-S	91.96	90.76	78.52	88.99	93.61	91.66	85.34	82.19	96.07	98.80	94.00	96.70	90.70	89.60
ConvNext-B	92.30	91.52	80.44	90.00	93.91	93.01	86.65	83.75	96.41	99.20	93.60	96.40	90.60	91.40
ViT-T	80.81	77.07	46.78	69.07	81.06	76.55	64.13	57.86	91.08	95.20	85.50	90.40	67.30	64.50
ViT-S	86.77	83.46	63.19	80.58	87.98	84.05	74.29	69.94	94.74	97.70	89.20	94.30	84.20	80.60
ViT-B	90.07	87.48	71.28	84.88	91.01	88.64	79.99	76.42	95.66	98.00	90.40	93.80	86.20	84.80
DeiT-T	80.68	77.21	51.24	71.83	80.30	76.83	65.03	59.54	89.94	91.60	86.50	90.50	73.90	73.20
DeiT-S	87.52	84.88	63.03	80.70	89.13	85.70	76.75	71.37	94.14	98.30	91.40	95.20	85.70	84.10
DeiT-B	89.66	86.79	68.77	84.26	91.10	89.19	80.31	77.25	95.38	98.80	92.30	96.00	86.70	84.30
VMamba-T(v0)	90.03	89.31	70.38	85.59	91.49	89.59	82.01	78.63	95.73	98.00	91.60	94.90	86.10	86.70
VMamba-S(v0)	90.53	90.76	73.39	87.78	93.13	90.71	83.92	80.42	96.16	99.10	92.80	95.40	89.40	88.50
VMamba-B(v0)	91.75	90.62	73.90	88.33	93.29	91.19	83.96	81.66	96.00	98.80	93.50	96.20	89.70	88.40
VMamba-T(v2)	91.15	89.87	75.18	87.41	92.09	91.06	83.66	79.71	95.84	98.50	92.20	96.30	87.20	86.80
VMamba-S(v2)	92.03	90.79	76.15	88.81	93.22	92.25	85.57	81.89	96.37	99.20	94.10	97.40	90.90	89.50
VMamba-B(v2)	92.37	91.27	77.30	89.11	93.70	92.64	86.03	83.62	96.37	99.10	94.00	96.50	90.80	89.80
MambaVision-T	90.67	88.83	73.07	86.40	91.93	90.07	81.87	78.42	95.73	98.60	93.70	96.60	89.10	87.70
MambaVision-S	91.22	90.19	75.18	88.19	92.78	91.19	84.01	81.18	96.03	99.40	94.40	97.80	91.40	90.10
MambaVision-B	91.77	90.65	78.42	89.18	93.70	92.81	86.35	83.78	96.30	99.10	94.60	97.20	91.40	90.60
Swin-T	90.05	88.83	71.51	86.19	91.08	88.94	79.39	76.49	95.27	97.90	91.70	95.30	85.50	84.00
Swin-S	90.67	88.86	73.35	87.25	91.91	89.68	81.55	78.81	96.25	98.30	91.80	95.50	86.10	85.40
Swin-B	91.08	89.96	75.09	87.87	92.62	91.22	83.43	80.65	95.95	98.60	92.30	95.60	89.20	87.40

Table 10. Average Precision (AP) scores for different architectures on the COCO-DC dataset [35], detailing results for small (APs), medium (APm), and large objects (API).

Model	AP	APs	APm	API	AP	APs	APm	API	AP	APs	APm	API	AP	APs	APm	API
	Original				Color				Texture				Average			
ConvNext-T	66.2	41.0	61.3	71.4	55.0 _(-11.20)	26.3 _(-14.70)	49.0 _(-12.30)	60.6 _(+10.80)	53.5 _(-12.70)	26.5 _(-14.50)	47.5 _(-13.80)	60.5 _(-10.90)	54.25 _(-11.95)	26.40 _(-14.60)	48.25 _(-13.05)	60.55 _(-10.85)
Swin-T	66.3	45.1	61.6	71.6	55.1 _(-11.20)	29.2 _(-15.90)	48.0 _(-13.60)	62.0 _(+09.60)	53.6 _(-12.70)	31.4 _(-13.70)	45.2 _(-16.40)	61.4 _(-10.20)	54.35 _(-11.95)	30.30 _(-14.80)	46.60 _(-15.00)	61.70 _(-09.90)
Swin-S	69.1	43.7	62.3	75.9	57.4 _(-11.70)	30.1 _(-13.60)	49.1 _(-13.20)	64.3 _(+11.60)	56.0 _(-13.10)	26.8 _(-16.90)	45.4 _(-16.90)	64.3 _(-11.60)	56.70 _(-12.40)	28.45 _(-15.25)	47.25 _(-15.05)	64.30 _(-11.60)
VSSM-T	69.3	47.1	64.4	75.5	56.2 _(-13.10)	30.9 _(-16.20)	48.4 _(-16.00)	63.7 _(+11.80)	53.2 _(-16.10)	29.2 _(-17.90)	44.9 _(-19.50)	61.0 _(-14.50)	54.57 _(-14.73)	30.05 _(-17.05)	46.65 _(-17.75)	62.35 _(-13.15)
VSSM-S	70.9	51.7	64.6	77.1	56.7 _(-14.20)	31.3 _(-20.40)	48.9 _(-15.70)	64.3 _(-12.80)	55.7 _(-15.20)	26.8 _(-24.90)	47.4 _(-17.20)	64.4 _(-12.70)	56.20 _(-14.70)	29.05 _(-22.65)	48.15 _(-16.45)	64.35 _(-12.75)

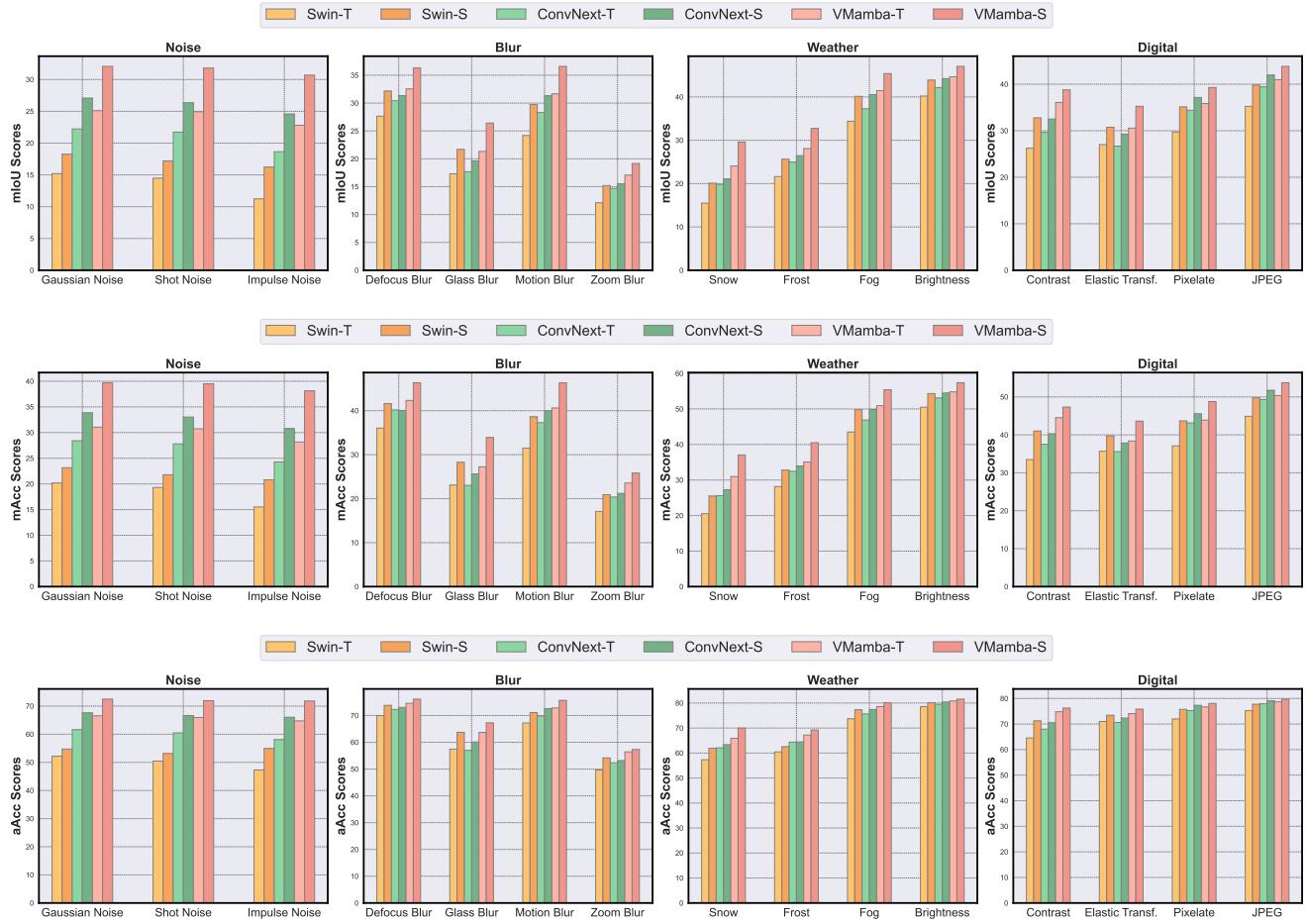


Figure 24. mIoU, mAcc, and aAcc score for different architectures on AED20k-C dataset

Table 11. Robust accuracy of various architectures against low-frequency and high-frequency-based perturbation using PGD attack at varying level of perturbation budget $\epsilon \in \{1/255, 2/255, 4/255, 8/255, 12/255, 16/255\}$.

Filter →	Low-Pass						High-Pass						All-Pass		
$\epsilon \rightarrow$	1/255	2/255	4/255	8/255	12/255	16/255	1/255	2/255	4/255	8/255	12/255	16/255	1/255	2/255	4/255
ResNet-50	94.94	93.28	91.38	88.42	87.74	87.86	72.68	43.58	28.70	16.38	11.62	10.76	29.68	2.04	0.08
ConvNext-T	95.98	95.46	94.60	91.76	89.70	89.50	58.26	26.46	13.62	4.40	2.38	2.06	20.18	1.60	0.02
ConvNext-S	96.52	96.14	95.80	93.48	91.70	91.36	66.58	34.66	20.50	7.38	3.66	3.12	30.10	3.98	0.08
ConvNext-B	96.84	96.76	96.14	93.62	92.40	91.72	64.96	39.26	22.98	8.76	5.46	4.56	31.02	6.24	0.18
ViT-T	85.52	77.56	68.50	53.02	47.28	46.32	69.32	38.48	19.24	5.52	2.56	1.88	11.14	0.16	0.00
ViT-S	93.12	88.60	83.50	72.84	69.22	69.06	77.14	46.48	27.94	11.08	6.76	5.28	18.74	0.42	0.00
ViT-B	94.66	90.28	87.40	80.68	77.50	77.26	85.74	61.84	46.52	26.18	18.28	16.28	30.94	1.30	0.08
VMamba-T	96.46	96.04	95.22	93.02	91.70	91.18	58.12	32.50	16.46	5.74	2.92	2.82	24.00	5.16	0.28
VMamba-S	97.02	96.90	96.48	94.10	92.62	92.24	67.64	44.00	28.18	11.36	6.46	5.48	32.54	10.90	0.84
VMamba-B	97.28	97.00	96.56	94.28	92.90	92.24	66.60	45.44	27.68	10.36	6.08	4.78	33.10	9.88	0.40
MambaVision-T	96.44	95.98	95.18	93.02	92.40	92.54	60.52	35.32	22.90	9.00	4.88	3.86	18.38	3.04	0.16
MambaVision-S	96.74	96.26	95.44	93.56	92.84	93.10	67.60	44.56	33.24	17.48	11.14	9.24	23.80	4.46	0.30
MambaVision-B	96.80	96.24	95.82	93.86	93.04	93.22	70.28	45.68	33.20	18.36	12.86	11.24	28.80	6.40	0.72
Swin-T	96.24	95.86	95.24	93.28	91.94	91.30	49.12	23.42	9.78	2.06	1.10	0.88	12.60	1.86	0.00
Swin-S	96.92	96.60	96.16	94.26	93.30	92.70	60.70	36.50	20.18	7.48	3.90	3.30	25.92	6.74	0.42
Swin-B	96.82	96.94	96.24	94.72	94.10	93.34	61.04	40.50	24.68	9.30	5.82	4.34	30.14	9.76	0.70

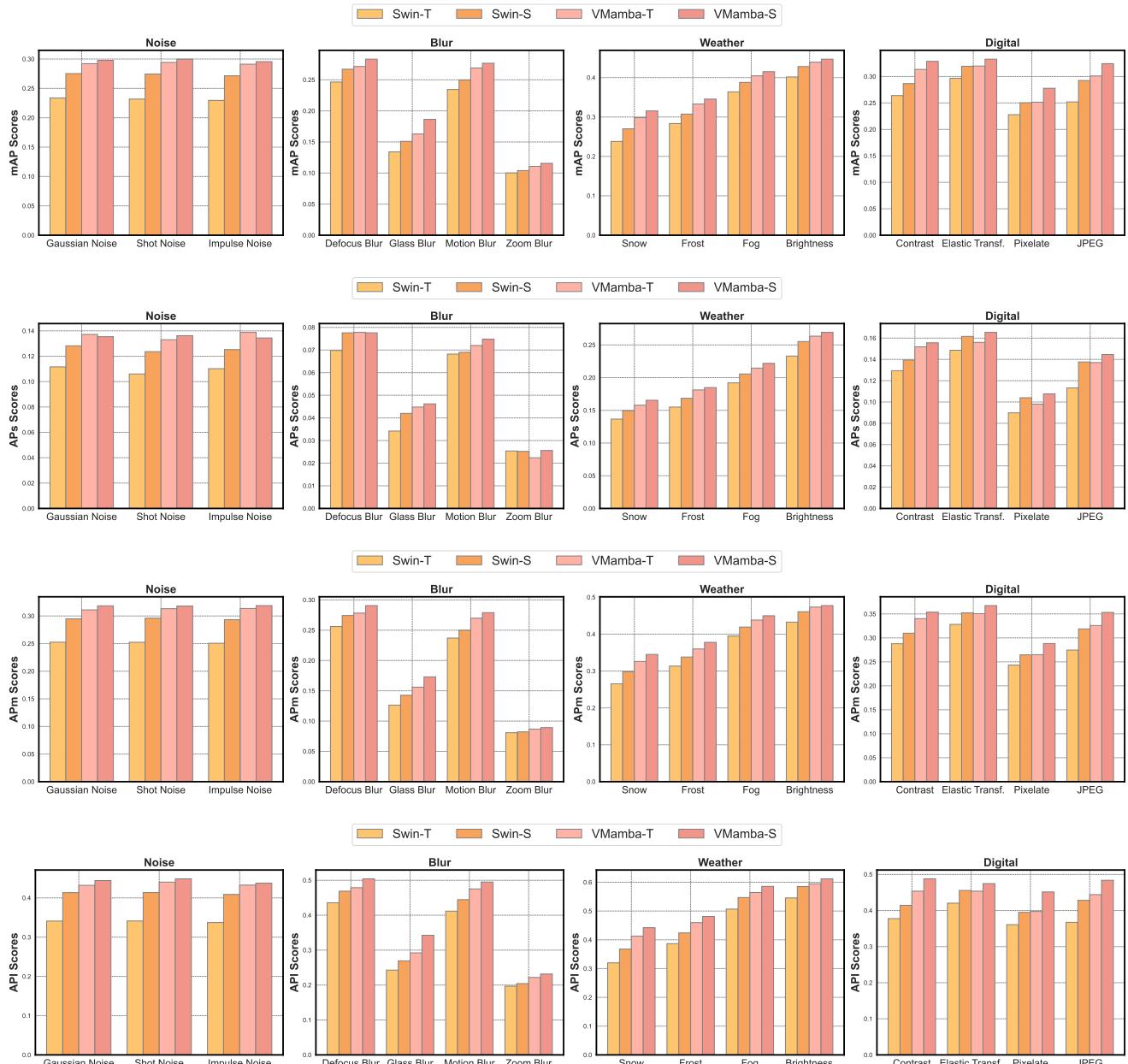


Figure 25. mAP , AP_s , AP_m , and AP_l score for different architectures on COCO-C dataset

Transferability Results for FGSM at Epsilon 8/255

	V-Mamba-T-v2	V-Mamba-S-v2	V-Mamba-B-v2	V-Mamba-T-v0	V-Mamba-S-v0	V-Mamba-B-v0	MambaVision-T	MambaVision-S	MambaVision-B	ResNet-18	ResNet-50	VGG-16	VGG-19	DenseNet-121	DenseNet-161	ConvNext-T	ConvNext-S	ConvNext-B	ViT-T	ViT-S	ViT-B	DeiT-T	DeiT-S	DeiT-B	Swin-T	Swin-S	Swin-B	
VMamba-T-v2	42.90	66.34	65.10	72.12	75.24	74.42	74.44	73.00	73.16	69.66	80.20	67.74	69.24	76.30	76.32	72.66	73.80	72.84	79.46	83.64	86.94	75.76	82.00	82.32	72.22	76.38	74.60	
VMamba-S-v2	62.24	48.42	63.00	69.40	71.34	70.78	71.34	71.92	71.34	68.66	79.70	66.28	67.44	75.26	75.14	71.40	71.04	70.18	78.42	81.56	84.58	75.96	80.30	80.22	70.96	72.92	71.14	
VMamba-B-v2	65.52	66.96	51.24	72.50	74.38	72.94	75.22	73.82	73.08	70.04	81.22	68.62	69.56	77.34	76.84	73.54	73.20	72.32	79.12	83.62	86.24	76.20	81.58	81.82	73.24	76.30	73.88	
VMamba-T-v0	77.74	78.94	79.34	55.68	73.88	73.30	79.98	81.40	81.72	73.60	84.70	74.40	75.12	79.64	80.78	78.04	78.78	78.74	80.62	85.28	88.08	78.22	83.60	84.26	75.08	77.28	77.24	
VMamba-S-v0	76.50	77.10	77.20	70.82	58.28	71.52	78.58	79.28	79.36	72.22	83.54	72.68	73.04	78.42	79.06	76.34	76.12	76.00	79.72	82.92	85.62	76.88	81.78	82.68	74.22	75.82	74.84	
VMamba-B-v0	76.66	76.96	76.46	71.48	72.56	59.18	80.10	80.10	79.56	72.42	83.84	72.90	73.96	79.90	79.90	77.04	77.14	76.30	80.32	84.74	86.38	78.48	82.88	82.68	74.52	76.34	75.30	
MambaVision-T	77.74	79.04	79.78	78.52	81.48	81.36	46.18	73.30	75.64	68.22	80.08	69.66	70.98	77.02	78.08	79.50	80.88	81.82	78.30	84.20	87.86	75.50	82.08	84.64	77.78	81.24	80.80	
MambaVision-S	77.86	79.56	79.24	79.88	82.70	81.90	74.96	53.42	73.84	69.60	83.60	70.68	70.98	78.42	79.50	80.40	81.14	81.26	79.90	85.64	88.52	76.96	83.42	84.80	78.50	82.40	81.88	
MambaVision-B	75.90	77.84	77.00	78.96	81.50	81.16	75.46	71.96	52.68	70.16	83.46	70.04	71.76	78.46	79.22	80.52	79.92	79.82	85.18	88.12	77.76	83.80	84.10	78.46	80.66	80.30		
ResNet-18	85.48	88.34	89.24	84.48	88.74	89.40	83.22	86.48	88.86	1.50	71.46	51.34	51.98	57.74	63.34	84.84	87.06	88.66	74.12	85.98	91.00	70.20	85.42	89.36	82.94	88.14	88.16	
ResNet-50	81.38	83.24	83.84	80.56	85.22	85.36	81.06	82.54	84.88	55.88	30.46	60.84	61.44	61.22	66.36	80.30	82.20	83.38	75.94	84.74	89.12	72.18	84.12	87.54	80.64	85.00	85.42	
VGG-16	83.80	86.80	87.42	85.72	89.74	89.82	84.90	86.18	88.42	58.10	80.84	2.92	37.80	73.40	74.42	83.82	86.24	87.34	80.42	89.42	92.98	76.10	89.16	91.92	83.52	88.77	88.48	
VGG-19	83.92	86.80	87.14	85.68	89.44	89.30	84.80	85.08	87.48	57.50	79.24	36.16	4.54	71.66	72.30	83.70	85.66	86.86	79.60	88.66	92.24	75.80	88.36	91.70	80.70	88.64	88.46	
DenseNet-121	85.42	87.82	88.80	84.68	88.10	88.84	83.74	86.18	88.48	54.10	71.70	62.36	63.40	15.36	64.58	84.06	86.14	88.00	76.38	85.96	90.72	72.52	86.28	89.74	83.34	88.32	88.72	
DenseNet-161	79.60	82.02	82.94	79.22	83.30	83.44	78.76	81.30	83.92	52.12	68.46	55.24	57.54	55.58	55.88	8.74	78.46	80.04	81.42	72.88	82.78	86.22	71.48	82.86	86.18	78.24	83.68	84.04
ConvNext-T	69.00	71.46	71.18	68.56	71.82	70.64	73.50	73.92	74.40	67.60	77.96	65.22	65.54	72.54	72.56	76.36	81.96	63.76	76.92	82.74	85.58	74.10	79.88	80.78	67.06	71.88	70.78	
ConvNext-S	69.54	70.62	70.48	69.90	71.66	70.74	74.68	74.38	75.24	70.70	79.48	68.44	69.78	74.68	74.16	63.48	49.10	63.62	78.48	82.94	85.02	75.70	80.36	81.16	69.24	71.54	69.74	
ConvNext-B	70.54	71.78	69.78	71.46	73.32	71.44	76.78	75.02	74.82	71.42	81.72	69.84	71.28	77.20	76.72	67.34	66.24	51.32	80.26	83.98	86.22	77.50	82.64	81.36	69.86	73.44	70.84	
ViT-T	85.92	88.30	88.88	81.76	84.92	85.84	82.16	85.06	87.24	64.88	82.68	68.24	70.78	74.24	76.88	85.08	86.56	88.06	2.28	50.04	69.90	39.14	52.44	66.72	75.66	79.72	82.32	
ViT-S	81.74	83.76	84.46	77.56	80.64	80.92	78.94	80.78	82.68	66.62	81.70	68.32	70.58	73.70	76.28	82.34	82.78	84.38	45.40	11.02	54.82	51.84	55.50	60.90	72.50	75.98	77.88	
ViT-B	82.46	84.02	84.48	79.30	81.48	80.80	80.40	81.82	82.46	69.00	81.90	69.34	71.82	75.10	76.44	82.04	82.24	83.62	58.50	53.84	24.24	62.14	64.80	65.64	75.82	77.62	78.54	
DeiT-T	83.34	86.40	87.74	77.82	82.88	83.74	79.18	83.14	86.66	61.36	80.06	64.82	67.28	69.86	75.44	82.06	84.40	86.68	38.14	54.82	73.58	10.96	42.92	59.50	70.86	77.84	80.68	
DeiT-S	80.68	82.72	82.22	75.96	78.62	79.28	77.18	80.16	81.08	68.68	81.28	70.58	71.98	74.94	77.34	80.18	80.78	82.46	55.94	61.74	71.36	51.70	31.60	58.74	71.56	74.16	76.02	
DeiT-B	81.74	83.42	83.28	79.56	81.86	80.90	80.40	81.66	82.34	70.62	83.50	71.34	74.14	77.46	79.28	81.92	83.00	83.42	65.16	70.62	75.38	61.40	64.84	42.16	76.52	78.42	78.84	
Swin-T	72.40	76.06	75.58	67.86	72.58	71.80	76.80	77.16	78.08	70.32	82.82	70.50	71.82	76.88	78.08	71.46	72.12	71.96	76.48	80.66	85.46	72.52	78.12	80.16	28.86	56.22	55.12	
Swin-S	78.24	79.10	79.30	73.96	76.84	75.86	80.12	81.90	81.40	74.24	85.46	75.34	76.30	81.46	82.32	77.96	77.10	77.02	79.38	83.04	86.58	77.40	81.64	82.34	61.94	48.00	63.90	
Swin-B	78.88	79.40	79.04	76.08	77.68	77.24	81.48	82.16	81.04	74.58	86.60	75.38	77.26	82.18	81.96	78.16	78.46	77.32	81.48	84.74	87.64	79.18	83.56	82.86	66.82	68.28	54.76	

Figure 26. Robust accuracy of various architectures under white-box and black-box settings for FGSM attack. Adversarial examples are crafted at a perturbation budget $\epsilon = \frac{8}{255}$.

Transferability Results for PGD at Epsilon 8/255

	V-Mamba-T-v2	V-Mamba-S-v2	V-Mamba-B-v2	V-Mamba-T-v0	V-Mamba-S-v0	V-Mamba-B-v0	MambaVision-T	MambaVision-S	MambaVision-B	ResNet-18	ResNet-50	VGG-16	VGG-19	DenseNet-121	DenseNet-161	ConvNext-T	ConvNext-S	ConvNext-B	ViT-T	ViT-S	ViT-B	DeiT-T	DeiT-S	DeiT-B	Swin-T	Swin-S	Swin-B
VMamba-T-v2	0.01	37.26	35.56	64.48	73.60	70.60	70.86	67.72	68.78	69.58	81.44	54.40	54.40	75.30	77.72	63.54	67.18	65.83	82.58	89.16	91.80	78.82	87.98	89.02	65.76	74.72	71.70
VMamba-S-v2	16.88	0.02	23.94	53.54	60.70	58.56	60.98	60.22	59.16	65.14	77.54	45.10	48.76	70.98	72.44	56.24	57.72	56.16	79.26	86.08	89.42	76.46	85.26	86.02	58.92	65.92	61.46
VMamba-B-v2	15.98	25.26	0.00	54.46	63.08	58.60	64.56	59.68	59.34	65.72	77.38	44.18	48.36	70.88	72.30	56.30	57.94	53.42	79.76	86.24	89.94	75.68	85.08	85.96	57.40	66.52	61.74
VMamba-T-v0	71.14	75.74	75.30	72.72	53.94	51.38	78.72	81.60	82.90	71.56	83.84	65.42	67.22	76.94	80.16	69.12	71.16	72.26	81.16	87.74	91.34	77.44	87.12	88.90	58.96	67.66	67.34
VMamba-S-v0	63.50	66.88	66.50	27.90	0.32	37.58	72.32	75.14	76.78	66.84	79.44	58.16	60.44	72.32	75.74	59.74	58.68	60.84	78.30	84.62	89.16	74.66	84.10	86.04	49.70	56.82	55.90
VMamba-B-v0	61.42	65.24	63.00	25.50	37.36	0.40	73.98	75.74	77.20	68.64	80.96	58.46	61.06	73.90	76.92	59.60	61.08	59.96	79.16	85.18	89.48	74.82	84.36	85.52	46.96	55.02	53.64
MambaVision-T	76.62	79.98	80.92	79.74	83.68	83.86	0.00	64.18	71.78	71.42	85.76	66.84	69.54	79.00	81.82	80.42	82.50	84.24	82.62	88.90	92.46	78.78	87.60	89.40	78.86	83.42	83.14
MambaVision-S	67.52	73.60	71.86	76.54	81.78	81.22	55.92	0.00	53.36	71.32	83.72	61.52	63.54	77.44													

Transferability Results for MIFGSM at Epsilon 8/255

	VSSM-T_v2	VSSM-S_v2	VSSM-B_v2	VSSM-T_v0	VSSM-S_v0	VSSM-B_v0	ResNet-18	ResNet-50	VGG-16	VGG-19	DenseNet-121	DenseNet-161	ConvNext-T	ConvNext-S	ConvNext-B	ViT-T	ViT-S	DeiT-T	DeiT-S	DeiT-B	Swin-T	Swin-S	Swin-B	
VSSM-T_v2	0.00	25.98	24.92	52.08	60.62	58.12	61.40	71.54	46.10	49.04	67.08	67.76	50.82	53.96	51.98	76.42	81.64	86.56	71.92	80.82	81.12	53.78	64.04	60.42
VSSM-S_v2	12.46	0.00	16.98	40.32	46.72	44.42	55.16	65.62	36.18	40.58	60.34	58.90	42.32	43.14	42.28	71.84	76.92	82.58	67.78	75.24	76.08	46.86	53.04	49.54
VSSM-B_v2	12.90	18.50	0.00	42.72	50.88	46.12	56.54	67.94	36.50	41.04	62.00	61.26	43.82	45.70	41.18	72.82	78.54	83.54	68.50	77.04	77.28	47.12	55.76	50.32
VSSM-T_v0	59.86	63.26	63.04	0.00	44.14	42.98	64.84	76.82	60.32	62.28	69.80	72.48	57.82	58.60	60.20	74.98	81.18	85.80	71.04	78.46	80.74	50.80	57.62	56.68
VSSM-S_v0	50.74	52.68	53.06	22.84	0.04	26.14	60.38	71.66	52.62	53.94	63.54	65.84	48.94	45.14	47.14	70.62	75.90	82.38	66.60	73.42	75.66	39.82	46.14	44.72
VSSM-B_v0	48.90	50.46	48.94	20.80	26.24	0.02	60.30	72.68	49.90	52.60	65.86	66.90	46.60	47.72	46.90	71.42	77.90	82.80	68.04	75.16	76.24	37.74	45.26	43.44
ResNet-18	82.26	85.36	86.54	81.90	86.70	87.60	0.00	58.10	34.42	36.24	38.94	46.58	80.26	83.98	86.60	71.20	84.92	90.38	68.20	84.80	89.28	79.56	86.46	87.14
ResNet-50	73.02	76.76	77.84	75.32	80.22	80.78	42.86	0.02	45.18	46.28	38.78	48.32	72.02	75.62	77.76	73.32	82.56	88.16	68.78	83.22	86.44	75.32	81.96	82.26
VGG-16	78.94	83.72	83.52	83.52	88.52	88.38	55.24	77.46	0.10	8.74	71.58	70.18	78.68	82.04	83.12	80.52	89.44	93.04	75.90	89.58	92.78	80.84	87.66	87.48
VGG-19	78.28	82.64	82.98	82.16	87.26	87.30	51.86	74.32	6.10	0.08	67.12	66.78	77.54	81.24	82.52	79.88	89.10	92.48	74.70	89.02	91.44	80.30	87.10	87.28
DenseNet-121	80.42	84.10	84.72	80.20	84.52	85.54	38.82	56.54	49.74	51.16	0.00	44.26	78.58	81.78	84.36	73.24	83.98	89.80	69.74	84.50	87.92	78.42	85.90	85.66
DenseNet-161	65.66	70.14	71.66	67.02	73.34	75.16	35.96	46.92	35.74	37.70	27.28	0.20	63.50	66.68	68.54	67.84	78.42	84.62	65.86	78.14	82.22	68.06	75.68	76.02
ConvNext-T	47.00	52.92	52.02	45.18	53.72	51.30	58.44	69.50	45.16	47.28	61.18	61.00	0.00	22.54	29.56	72.84	79.22	84.36	68.52	76.80	78.06	44.20	55.56	53.36
ConvNext-S	43.58	47.44	46.90	40.28	45.20	44.92	56.94	66.66	43.48	46.30	57.70	58.82	17.96	0.00	21.86	70.16	76.50	81.10	67.06	74.86	75.44	39.50	48.78	45.42
ConvNext-B	46.04	50.72	47.28	45.90	52.08	48.42	62.32	70.76	47.62	51.42	65.76	64.62	26.30	25.46	0.00	74.56	79.92	84.26	71.48	78.76	78.28	42.04	52.54	47.72
ViT-T	88.00	90.46	91.30	83.88	87.18	88.38	63.22	82.38	65.34	67.30	73.98	77.58	87.32	88.62	90.92	0.00	44.78	68.80	31.50	48.70	67.86	78.08	82.98	85.36
ViT-S	83.24	84.96	86.11	77.26	81.40	81.92	63.94	79.90	65.38	67.26	72.66	74.94	82.24	83.62	85.40	31.34	0.00	39.28	43.48	44.60	54.76	70.90	74.54	77.28
ViT-B	81.56	83.40	84.32	75.80	78.74	79.74	64.66	79.78	64.98	67.86	71.40	74.14	80.88	81.98	83.42	44.58	31.26	0.00	52.74	54.14	54.22	72.16	74.86	76.54
DeiT-T	85.78	89.00	90.22	80.72	84.56	86.72	58.76	79.88	61.58	63.88	68.38	75.84	85.32	87.22	89.28	25.02	53.76	73.02	0.00	30.00	56.46	73.34	81.12	84.16
DeiT-S	78.76	81.52	82.96	72.66	76.38	77.54	60.46	77.64	62.42	65.02	68.58	73.56	78.60	80.36	82.32	31.04	42.84	62.16	21.38	0.00	34.40	64.58	70.02	73.38
DeiT-B	77.34	79.82	80.36	70.80	74.82	75.40	61.90	77.00	62.38	64.76	69.22	73.30	77.98	78.90	80.54	46.14	52.56	61.28	38.70	32.04	0.00	64.68	68.80	70.26
Swin-T	69.68	76.72	75.62	60.42	69.62	67.12	69.06	83.88	66.18	69.14	76.22	78.98	65.84	67.56	67.26	76.08	82.54	88.16	72.46	80.66	83.20	0.00	26.94	24.22
Swin-S	73.28	77.02	76.42	61.70	68.58	67.28	70.64	84.24	67.88	69.86	78.54	79.98	69.88	68.82	69.18	77.08	82.16	87.36	73.88	80.94	82.48	18.92	0.00	21.70
Swin-B	72.88	75.82	74.74	64.26	70.88	68.36	71.16	85.42	66.70	69.92	79.02	80.38	69.36	70.58	68.00	79.40	84.22	88.66	76.00	83.24	84.26	20.90	27.08	0.00

Figure 28. Robust accuracy of various architectures under white-box and black-box settings for MIFGSM attack. Adversarial examples are crafted at a perturbation budget $\epsilon = \frac{8}{255}$.