

---

# “Rubik’s Cube: High-Order Channel Interactions with a Hierarchical Receptive Field” Supplementary Material

---

Anonymous Author(s)

Affiliation

Address

email

- 1 This supplementary document is organized as follows:  
2 Section 1 illustrates the process of how to form a hierarchical receptive field within the combination  
3 of the shifting and interaction operations.  
4 Section 2 provides the implementation details of Rubik’s cube convolution within the image restora-  
5 tion baselines.  
6 Section 3 provides the evaluation of our proposed Rubik’s cube convolution on the classification task.  
7 We conduct the experiments on three widely-used classification benchmarks with three representative  
8 baselines.  
9 Section 4 provides more quantitative and qualitative results.

## 10 1 The Hierarchical Receptive Field

11 As illustrated in Figure 1, given an input feature map  $\mathbf{X} \in \mathbb{R}^{H \times W \times C}$ , we evenly divide  $\mathbf{X}$  into five  
12 parts by the channel dimension, where the first is kept unchanged and the remaining four ones are  
13 shifted in a distinct spatial direction: left, right, top, and down. Subsequent to the shifting operation,  
14 we discard out-of-focus pixels and any vacant pixels are filled with zeros. The shifted feature  $\hat{\mathbf{X}}$  can  
15 be written as:

$$\begin{aligned}\hat{\mathbf{X}}[0 : H, 0 : W, 0 : C_{\text{id}}] &\leftarrow \mathbf{X}[0 : H, 0 : W, 0 : C_{\text{id}}], \\ \hat{\mathbf{X}}[0 : H, 1 : W, C_{\text{id}} : C_{\text{id}} + C_{\text{g}}] &\leftarrow \mathbf{X}[0 : H, 0 : W - 1, C_{\text{id}} : C_{\text{id}} + C_{\text{g}}], \\ \hat{\mathbf{X}}[0 : H, 0 : W - 1, C_{\text{id}} + C_{\text{g}} : C_{\text{id}} + 2C_{\text{g}}] &\leftarrow \mathbf{X}[0 : H, 1 : W, C_{\text{id}} + C_{\text{g}} : C_{\text{id}} + 2C_{\text{g}}], \\ \hat{\mathbf{X}}[0 : H - 1, 0 : W, C_{\text{id}} + 2C_{\text{g}} : C_{\text{id}} + 3C_{\text{g}}] &\leftarrow \mathbf{X}[1 : H, 0 : W, C_{\text{id}} + 2C_{\text{g}} : C_{\text{id}} + 3C_{\text{g}}], \\ \hat{\mathbf{X}}[1 : H, 0 : W, C_{\text{id}} + 3C_{\text{g}} : C_{\text{id}} + 4C_{\text{g}}] &\leftarrow \mathbf{X}[0 : H - 1, 0 : W, C_{\text{id}} + 3C_{\text{g}} : C_{\text{id}} + 4C_{\text{g}}],\end{aligned}\tag{1}$$

16 where  $C_{\text{id}}$  is the number of channels of the unchanged identity part,  $C_{\text{g}}$  is the number of channels of  
17 a shifted group, and  $C_{\text{id}} + 4 * C_{\text{g}} = C$ . Next, the shifted feature  $\hat{\mathbf{X}}$  is split into  $\hat{\mathbf{X}}_{\text{ori}} \in \mathbb{R}^{H \times W \times C_{\text{id}}}$   
18 and  $\{\hat{\mathbf{X}}_{\text{c1}}, \hat{\mathbf{X}}_{\text{c2}}, \hat{\mathbf{X}}_{\text{c3}}, \hat{\mathbf{X}}_{\text{c4}}\} \in \mathbb{R}^{H \times W \times C_{\text{g}}}$  along the channel dimension.

19 After the shifting operation, we construct the high-order interaction in the channel dimension as  
20 described in the manuscript. As shown in Figure 1, for the  $(i, j)$  pixel in the up-shifting group,  $\hat{\mathbf{X}}_{\text{c1}}$ ,  
21 it will interact with the  $(i + p, j)$  pixel in the down-shifting group,  $\hat{\mathbf{X}}_{\text{c12}}$ , and  $p$  indicates the number  
22 of the shifted pixels. Therefore, the combination of the shifting and interaction leads to a stratified  
23 receptive field along the channel dimension, reminiscent of a Rubik’s cube.

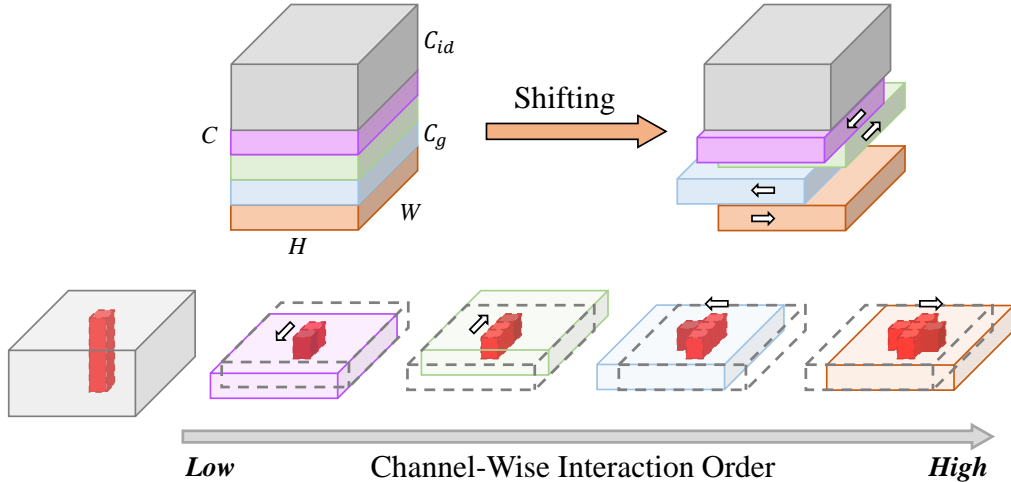


Figure 1: An illustration of the shifting operation and the formulation of the hierarchical receptive field. Specifically, the input feature is separated into five groups, where the last four are shifted into four direction and the first is unchanged. When a up-shifting group interacts the next down-shifting group, the  $(i, j)$  pixel will interweave with its neighboring pixel,  $(i + p, j)$ , where  $p$  denotes the number of shifted pixels. Therefore, with the combined action of the shifting and interaction operation, the receptive field will be expanded along the downward direction. The red-shaded region indicates the receptive field and the scarlet-shaded region presents the newly expanded receptive field after the corresponding interaction.

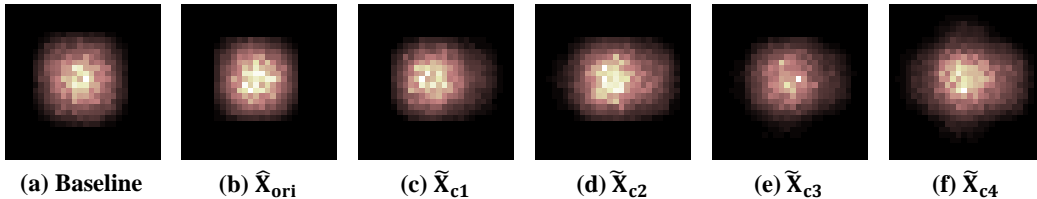


Figure 2: **Visualization of the effective receptive field [1].** (a) The baseline indicates the effective receptive field of the default architecture, and the last five describe the effective receptive field of the (b) identity group,  $\hat{X}_{ori}$ , (c) left-shifting group,  $\tilde{X}_{c1}$ , (d) right-shifting group,  $\tilde{X}_{c2}$ , (e) up-shifting group,  $\tilde{X}_{c3}$ , (f) and down-shifting group,  $\tilde{X}_{c4}$  after interactions in the Rubik’s cube convolution.

24 We visualize the effective receptive field [1] of the distinct groups in the Rubik’s cube convolution.  
 25 Figure 2 demonstrates the effective receptive field expands along the up, down, left, and right  
 26 directions after the corresponding shifting and interaction.

## 27 2 Implementation Details

28 Based on the competitive baselines, we create several variants of the baselines by replacing the  
 29 standard convolution with the proposed Rubik’s cube convolution:

- 30 1) **Original**: the baseline without any changes;
- 31 2) **RubikConv**: replacing the standard convolution in the original model with our designed  
 32 Rubik’s cube convolution;
- 33 3) **Conv1x1**: a baseline that replaces the RubikConv in the setting of 2) with four convolution  
 34 layers with  $1 \times 1$  kernel for a fair comparison with approximately the same number of  
 35 trainable parameters as 2).

36 Taking a network consisted of a stack of convolution layers with a  $3 \times 3$  kernel for example, we  
 37 replace the standard convolution layer with our proposed Rubik’s cube convolution layer in 2) or four  
 38 convolution with  $1 \times 1$  kernel in 3) from top to bottom.

### 39 **3 Evaluation on Image Classification**

40 Due to the limited space in the main manuscript, we provide the evaluation on the image clas-  
 41 sification task in the supplementary material. We choose three widely-used image classification  
 42 benchmarks: CIFAR-10 [2], CIFAR-100 [2], and CUB [3]. For comparison, we employ three  
 43 algorithms: AlexNet [4], VGG [5], and ResNet [6].

44 To validate the effectiveness of the proposed approach, we conduct extensive experiments as described  
 45 in the implementation details. The quantitative results are presented in Table 1, where the best results  
 46 are highlighted in bold. For the three recognition models across three benchmarks, integrating our  
 47 Rubik’s cube convolution operation into the baseline will achieve the performance improvement,  
 48 demonstrating the effectiveness of our operation in the recognition task.

Table 1: Quantitative comparison of image classification. “Acc-1” and “Acc-5” indicate the top-1 and top-5 classification accuracy.

Model	Metric	CIFAR-10			CIFAR-100			CUB		
		Original	Conv1x1	RubikConv	Original	Conv1x1	RubikConv	Original	Conv1x1	RubikConv
AlexNet	Acc-1	77.98	75.48	<b>82.63</b>	49.98	48.27	<b>52.77</b>	61.99	60.05	<b>64.31</b>
	Acc-5	98.69	95.71	<b>98.94</b>	70.43	66.41	<b>76.12</b>	83.48	83.09	<b>86.73</b>
VGG-16	Acc-1	84.62	83.54	<b>87.40</b>	55.76	55.06	<b>57.39</b>	79.70	78.52	<b>80.79</b>
	Acc-5	99.22	96.09	<b>99.26</b>	76.36	76.43	<b>79.20</b>	88.34	86.85	<b>89.06</b>
ResNet-18	Acc-1	89.45	87.85	<b>91.90</b>	59.63	58.74	<b>62.06</b>	85.72	84.16	<b>87.05</b>
	Acc-5	99.65	96.83	<b>99.70</b>	82.08	80.47	<b>83.90</b>	93.60	92.87	<b>94.65</b>

### 49 **4 Experiments**

50 **Quantitative Comparison.** Due to the limited space, we present the comparison on the World-III  
 51 dataset in the supplementary material. As the experiments on the manuscript, we adopt three baselines  
 52 (PanNet [7], MulNet [8], and INNFormer [9]) for evaluation. We conduct experiments as described  
 53 in the implementation details. As described in Table 2, we observe a performance gain by integrating  
 54 our proposed Rubik’s cube convolution across all competitive baselines.

Table 2: Quantitative comparisons of pan-sharpening.

Model	Configurations	WorldView-III			
		PSNR $\uparrow$	SSIM $\uparrow$	SAM $\downarrow$	ERGAS $\downarrow$
PanNet	Original	29.6863	0.9072	0.0853	3.4260
	Conv1x1	29.4305	0.8973	0.1008	3.6954
	RubikConv	<b>39.9831</b>	<b>0.9139</b>	<b>0.0812</b>	<b>3.2453</b>
MulNet	Original	30.4807	0.9211	0.0769	3.1196
	Conv1x1	30.4186	0.9131	0.0837	3.3364
	RubikConv	<b>30.6426</b>	<b>0.9231</b>	<b>0.0754</b>	<b>3.0835</b>
INNFormer	Original	30.4349	0.9204	0.0756	3.1439
	Conv1x1	30.3850	0.9175	0.0827	3.3061
	RubikConv	<b>30.5052</b>	<b>0.9211</b>	<b>0.0742</b>	<b>3.1118</b>

55 **Qualitative Comparison.** Due to the limited space in the manuscript, we present more visualizations  
 56 in the supplementary materials. As illustrated in Figure 3 and 4, integrating our Rubik’s cube  
 57 convolution into baselines generate the visually pleasant enhanced results. Specifically, the baseline  
 58 and the baseline with Conv1x1 fails to recover texture and suffers from artifacts and color distortion.  
 59 In contrast, the baseline combined with our Rubik’s cube convolution operator achieves details  
 60 reconstruction, artifact reduction, and color consistency.

61 We also provide the visual comparison on the image de-noising task in Figure 5 and 6. The qualitative  
 62 results consistently demonstrate that the improvement on the visual quality by integrating our Rubik's  
 63 cube convolution. The comprehensive visual results in the supplementary material demonstrate the  
 64 effectiveness of our proposed operator.



Figure 3: Visual comparison of DRBN [10] on the LOL [11] dataset.

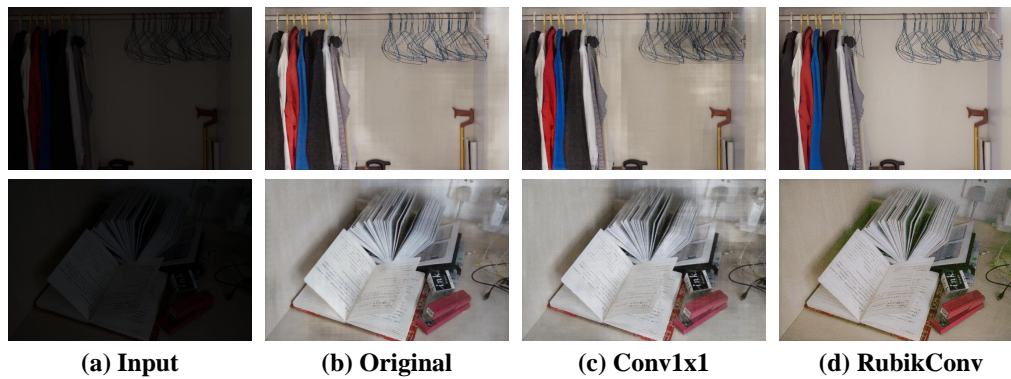


Figure 4: Visual comparison of SID [12] on the LOL [11] dataset.

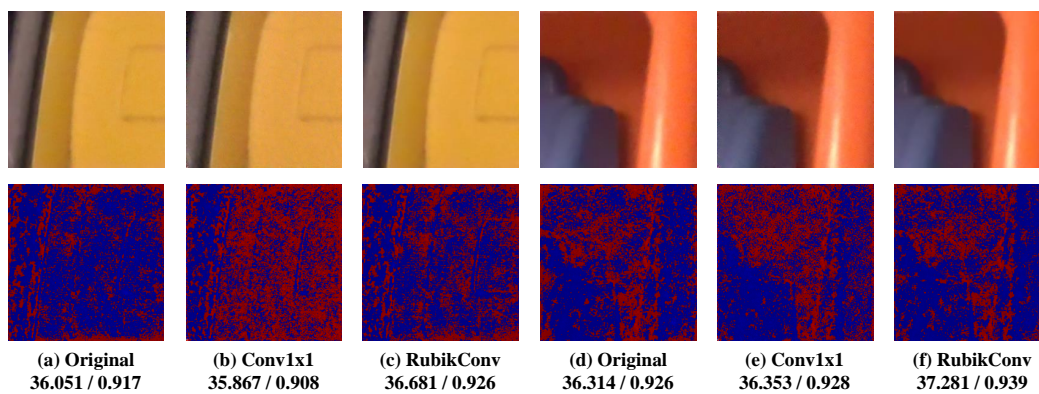


Figure 5: Visual comparison of DnCNN [13] on the SIDD [14] dataset. The second row presents the error map between the corresponding denoised results and the clean images. The numbers indicate the corresponding PSNR/SSIM metrics.

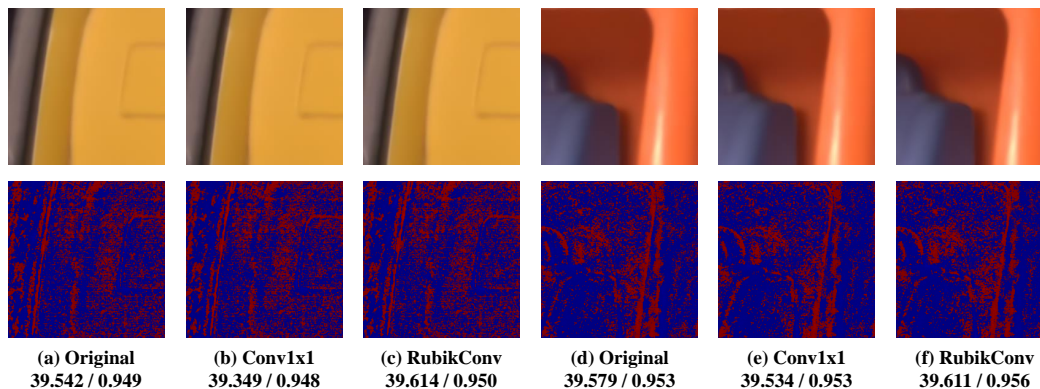


Figure 6: Visual comparison of MPRNet [15] on the SIDD [14] dataset. The second row presents the error map between the corresponding denoised results and the clean images. The numbers indicate the corresponding PSNR/SSIM metrics.

## 65 References

- 66 [1] Wenjie Luo, Yujia Li, Raquel Urtasun, and Richard Zemel. Understanding the effective receptive  
67 field in deep convolutional neural networks. *Advances in Neural Information Processing Systems*,  
68 29, 2016.
- 69 [2] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.  
70 2009.
- 71 [3] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The  
72 caltech-ucsd birds-200-2011 dataset. 2011.
- 73 [4] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep  
74 convolutional neural networks. *Communications of the ACM*, 60(6):84–90, 2017.
- 75 [5] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale  
76 image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- 77 [6] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
78 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
79 pages 770–778, 2016.
- 80 [7] Junfeng Yang, Xueyang Fu, Yuwen Hu, Yue Huang, Xinghao Ding, and John Paisley. PanNet:  
81 A deep network architecture for pan-sharpening. In *Proceedings of the IEEE International  
82 Conference on Computer Vision*, pages 5449–5457, 2017.
- 83 [8] Man Zhou, Keyu Yan, Jie Huang, Ziheng Yang, Xueyang Fu, and Feng Zhao. Mutual information-  
84 driven pan-sharpening. In *Proceedings of the IEEE/CVF Conference on Computer Vision and  
85 Pattern Recognition*, pages 1798–1808, 2022.
- 86 [9] Man Zhou, Jie Huang, Yanchi Fang, Xueyang Fu, and Aiping Liu. Pan-sharpening with  
87 customized transformer and invertible neural network. In *Proceedings of the AAAI Conference  
88 on Artificial Intelligence*, volume 36, pages 3553–3561, 2022.
- 89 [10] Wenhan Yang, Shiqi Wang, Yuming Fang, Yue Wang, and Jiaying Liu. From fidelity to percep-  
90 tual quality: A semi-supervised approach for low-light image enhancement. In *Proceedings  
91 of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages  
92 3063–3072, 2020.
- 93 [11] Chen Wei, Wenjing Wang, Wenhan Yang, and Jiaying Liu. Deep Retinex decomposition for  
94 low-light enhancement. *arXiv preprint arXiv:1808.04560*, 2018.
- 95 [12] Lijun Zhang, Xiao Liu, Erik Learned-Miller, and Hui Guan. Sid-nism: A self-supervised  
96 low-light image enhancement framework. *arXiv preprint arXiv:2012.08707*, 2020.

- 97 [13] Kai Zhang, Wangmeng Zuo, Yunjin Chen, Deyu Meng, and Lei Zhang. Beyond a gaussian  
98 denoiser: Residual learning of deep cnn for image denoising. *IEEE Transactions on Image*  
99 *Processing*, 26(7):3142–3155, 2017.
- 100 [14] Abdelrahman Abdelhamed, Stephen Lin, and Michael S Brown. A high-quality denoising  
101 dataset for smartphone cameras. In *Proceedings of the IEEE Conference on Computer Vision*  
102 *and Pattern Recognition*, pages 1692–1700, 2018.
- 103 [15] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, Ming-  
104 Hsuan Yang, and Ling Shao. Multi-stage progressive image restoration. In *Proceedings of*  
105 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14821–14831,  
106 2021.