# MIPI 2023 Challenge on RGBW Fusion: Methods and Results

Qianhui Sun    Qingyu Yang    Chongyi Li    Shangchen Zhou    Ruicheng Feng
Yuekun Dai    Wenxiu Sun    Qingpeng Zhu    Chen Change Loy    Jinwei Gu
Hongyuan Yu    Yuqing Liu    Weichen Yu    Lin Ge    Xiaolin Zhang    Qi Jia
Heng Zhang    Xuanwu Yin    Kunlong Zuo    Qi Wu    Wenjie Lin    Ting Jiang
Chengzhi Jiang    Mingyan Han    Xinpeng Li    Jinting Luo    Lei Yu    Haoqiang Fan
Shuaicheng Liu    Kunyu Wang    Chengzhi Cao    Yuanshen Guan    Jiyuan Xia
Ruikang Xu    Mingde Yao    Zhiwei Xiong

## Abstract

*Developing and integrating advanced image sensors with novel algorithms in camera systems are prevalent with the increasing demand for computational photography and imaging on mobile platforms. However, the lack of high-quality data for research and the rare opportunity for an in-depth exchange of views from industry and academia constrain the development of mobile intelligent photography and imaging (MIPI). With the success of the 1st MIPI Workshop@ECCV 2022, we introduce the second MIPI challenge, including four tracks focusing on novel image sensors and imaging algorithms. This paper summarizes and reviews the RGBW Joint Fusion and Denoise track on MIPI 2023. In total, 69 participants were successfully registered, and 4 teams submitted results in the final testing phase. The final results are evaluated using objective metrics, including PSNR, SSIM, LPIPS, and KLD. A detailed description of the models developed in this challenge is provided in this paper. More details of this challenge and the link to the dataset can be found at https://mipi-challenge.org/MIPI2023/.*

## 1. Introduction

RGBW is a new type of CFA (Color Filter Array) pattern (Fig. 1 (a)) designed for image quality enhancement under low light conditions. Thanks to the higher optical transmittance of white pixels over conventional red, green, and blue pixels, the signal-to-noise ratio (SNR) of images captured by this type of sensor increases significantly, thus boosting the image quality, especially under low light conditions. Recently, several phone OEMs [1, 2, 3] have adopted RGBW

Qianhui Sun[1] (sunqianhui@sensebrain.site), Qingyu Yang[1] (yangqingyu@sensebrain.site), Chongyi Li[4], Shangchen Zhou[4], Ruicheng Feng[4], Wenxiu Sun[2,3], Qingpeng Zhu[2], Chen Change Loy[4], Jinwei Gu[1,3] are the MIPI 2023 challenge organizers ([1] SenseBrain, [2] SenseTime Research and Tetras.AI, [3] Shanghai AI Laboratory, [4] Nanyang Technological University). The other authors participated in the challenge. Please refer to Appendix A for details.
MIPI 2023 challenge website: https://mipi-challenge.org/MIPI2023/

sensors in their flagship smartphones to improve the camera image quality.

The binning mode of RGBW is mainly used in the camera preview mode and video mode, in which a half-resolution Bayer is generated from the RGBW image, where spatial resolution is traded off for faster response. In this mode, every two pixels of the same color within a $2 \times 2$ window of the RGBW are averaged in the diagonal direction, and a diagonal-binning-Bayer image (DbinB) and a diagonal-binning-white image (DbinC) are generated. A fusion algorithm is demanded to enhance details and reduce noise in the Bayer image with the help of the white image (Fig. 1 (b)). A good fusion algorithm should be able to fully take advantage of the SNR and resolution benefit of white pixels.

The RGBW fusion problem becomes more challenging when the input DbinB and DbinC become noisy, especially under low light conditions. A joint fusion and denoise task is thus in demand for real-world applications.



Figure 1. The RGBW Fusion task: (a) the RGBW CFA. (b) In the binning mode, DbinB and DbinC are obtained by diagonal averaging of pixels of the same color within a $2\times2$ window. The joint fusion and denoise algorithm takes DbinB and DbinC as input to get a high-quality Bayer.

In this challenge, we intend to fuse DbinB and DbinC in Fig. 1 (b) to denoise and improve the Bayer. The solution is not necessarily learning-based. However, we provide a high-quality dataset of binning-mode RGBW input

Figure 2. An ISP to visualize the output Bayer and to calculate the loss function.

(DbinB and DbinC) and the output Bayer pairs to facilitate learning-based methods development, including 100 scenes (70 scenes for training, 15 for validation, and 15 for testing). The dataset is similar to the one provided in the first MIPI challenge, while we replaced some similar scenes with new ones. We also provide a simple ISP for participants to get the RGB image results from Bayer for quality assessment. Fig. 2 shows the pipeline of the simple ISP. The participants are also allowed to use other public-domain datasets. The algorithm performance is evaluated and ranked using objective metrics: Peak Signal-to-Noise Ratio (PSNR), Structural Similarity Index (SSIM) [15], Learned Perceptual Image Patch Similarity (LPIPS) [22], and KL-divergence (KLD).

We hold this challenge in conjunction with the second MIPI Challenge which will be held on CVPR 2023. Similar to the first MIPI challenges [7, 14, 17, 18, 16], we are seeking algorithms that fully take advantage of the SNR and resolution benefit of white pixels to enhance the final Bayer image in the binning model. MIPI 2023 consists of four competition tracks:

- **RGB+ToF Depth Completion** uses sparse and noisy ToF depth measurements with RGB images to obtain a complete depth map.

- **RGBW Sensor Fusion** fuses Bayer data and a monochrome channel data into Bayer format to increase SNR and spatial resolution.

- **RGBW Sensor Remosaic** converts RGBW RAW data into Bayer format so that it can be processed by standard ISPs.

- **Nighttime Flare Removal** is to improve nighttime image quality by removing lens flare effects.

## 2. MIPI 2023 RGBW Sensor Fusion

To facilitate the development of high-quality RGBW fusion solutions, we provide the following resources for participants:

- A high-quality dataset of aligned RGBW (DbinB and DbinC in Fig. 1 (b)) and Bayer. We enriched the scenes compared to the first MIPI challenge dataset. As far as we know, this is the only dataset consisting of aligned RGBW and Bayer pairs;

- A script that reads the provided raw data to help participants get familiar with the dataset;

- A simple ISP including basic ISP blocks to visualize the algorithm outputs and to evaluate image quality on RGB results;

- A set of objective image quality metrics to measure the performance of a developed solution.

### 2.1. Problem Definition

The RGBW fusion task aims to fuse the DbinB and DbinC of RGBW (Fig. 1 (b)) to improve the image quality of the Bayer output. By incorporating the white pixels in DbinC of higher spatial resolution and higher SNR, the output Bayer would potentially have better image quality. In addition, the binning mode of RGBW is mainly used for the preview and video modes in smartphones, thus requiring the fusion algorithms to be lightweight and power-efficient. While we do not rank solutions based on the running time or memory footprint, the computational cost is one of the most important criteria in real applications.

### 2.2. Dataset: Tetras-RGBW-Fusion

The training data contains 70 scenes of aligned RGBW (DbinB and DbinC input) and Bayer (ground-truth) pairs. DbinB at 0dB is used as the ground truth for each scene. Noise is synthesized on the 0dB DbinB and DbinC data to provide the noisy input at 24dB and 42dB, respectively. The synthesized noise consists of read noise and shot noise, and the noise models are calibrated on an RGBW sensor. The data generation steps are shown in Fig. 3. The testing data includes DbinB and DbinC inputs of 15 scenes at 24dB and 42dB, and the ground truth Bayer results are hidden from participants during the testing phase.

### 2.3. Evaluation

The evaluation consists of (1) the comparison of the fusion output Bayer and the reference ground truth Bayer, and
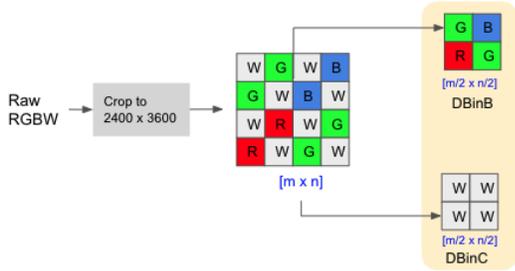
2

Figure 3. Data generation of the RGBW fusion task. The RGBW raw data is captured using an RGBW sensor and cropped into a size of 2400 × 3600. A Bayer (DbinB) and white (DbinC) image are obtained by averaging the same color in the diagonal direction within a 2 × 2 block.

(2) the comparison of RGB from the predicted and ground truth Bayer using a simple ISP (the code of the simple ISP is provided). We use

1. Peak Signal-to-Noise Ratio (PSNR)

2. Structural Similarity Index Measure (SSIM) [15]

3. Learned Perceptual Image Patch Similarity (LPIPS) [22]

4. Kullback–Leibler Divergence (KLD)

to evaluate the fusion performance. The PSNR, SSIM, and LPIPS will be applied to the RGB from the Bayer using the provided simple ISP code, while KLD is evaluated on the predicted Bayer directly.

A metric weighting PSNR, SSIM, KLD, and LPIPS is used to give the final ranking of each method, and we will report each metric separately as well. The code to calculate the metrics is provided. The weighted metric is shown below. The M4 score is between 0 and 100, and the higher score indicates the better overall image quality.

$$M4 = PSNR \cdot SSIM \cdot 2^{1-LPIPS-KLD}. \quad (1)$$

For each dataset, we report the average score over all the processed images belonging to it.

## 2.4. Challenge Phase

The challenge consisted of the following phases:

1. Development: The registered participants get access to the data and baseline code, and are able to train the models and evaluate their running time locally.

2. Validation: The participants can upload their models to the remote server to check the fidelity scores on the validation dataset, and to compare their results on the validation leaderboard.

3. Testing: The participants submit their final results, code, models, and factsheets.

## 3. Challenge Results

Four teams submitted their results in the final phase, which have been verified using their submitted code. Table. 1 summarizes the results in the final test phase. **RUSH MI**, **MegNR**, and **USTC-Zhalab** are the top three teams ranked by M4 are presented in Eq. (1), and **RUSH MI** shows the best overall performance. The proposed methods are described in Section 4, and the team members and affiliations are listed in Appendix A.

| Team name | PSNR | SSIM | LPIPS | KLD | M4 |
|---|---|---|---|---|---|
| RUSH MI | 38.587 | 0.977 | 0.0661 | 0.0718 | **68.58** |
| MegNR | 37.822 | 0.966 | 0.0815 | 0.0717 | **65.84** |
| USTC-Zhalab | 37.323 | 0.965 | 0.0854 | 0.0767 | **64.67** |
| VIDAR | 37.160 | 0.968 | 0.1023 | 0.0698 | 63.98 |

Table 1. MIPI 2023 Joint RGBW Fusion and Denoise challenge results and final rankings. PSNR, SSIM, LPIPS, and KLD are calculated between the submitted results from each team and the ground truth data. A weighted metric, M4, by Eq. (1) is used to rank the algorithm performance, and the top three teams with the highest M4 are highlighted.

To learn more about the algorithm performance, we evaluated the qualitative image quality in Fig. 4 and Fig. 5 in addition to the objective IQ metrics. While all teams in Table 1 have achieved high PSNR and SSIM, detail loss can be found on the texts of the card in Fig. 4 and detail loss or false color can be found on the mesh of the chair in Fig. 5. When the input has a large amount of noise, oversmoothing tends to yield higher PSNR at the cost of detail loss perceptually.

| Team name | 1200×1800 (measured) | 16M (estimated) |
|---|---|---|
| RUSH MI | **0.45s** | **3.33s** |
| MegNR | 8.46s | 62.60s |
| USTC-Zhalab | 69.60s | 515.04s |

Table 2. Running time of the top three solutions ranked by Eq. (1) in the MIPI 2023 Joint RGBW Fusion and Denoise challenge. The running time of input of 1200 × 1800 was measured, while the running time of a 64M RGBW sensor was based on estimation (the binning-mode resolution of a 64M RGBW sensor is 16M). The measurement was taken on an NVIDIA Tesla V100-SXM2-32GB GPU.

In addition to benchmarking the image quality of fusion algorithms, computational efficiency is evaluated because of the wide adoption of RGBW sensors in smartphones. We measured the running time of the RGBW fusion solutions of the top three teams in Table 2. While running time is not employed in the challenge to rank fusion algorithms, the computational cost is critical when developing smartphone algorithms. RUSH MI achieved the shortest running time among the top three solutions on a workstation GPU (NVIDIA Tesla V100-SXM2-32GB). With the sensor resolution of mainstream smartphones reaching 64M or even
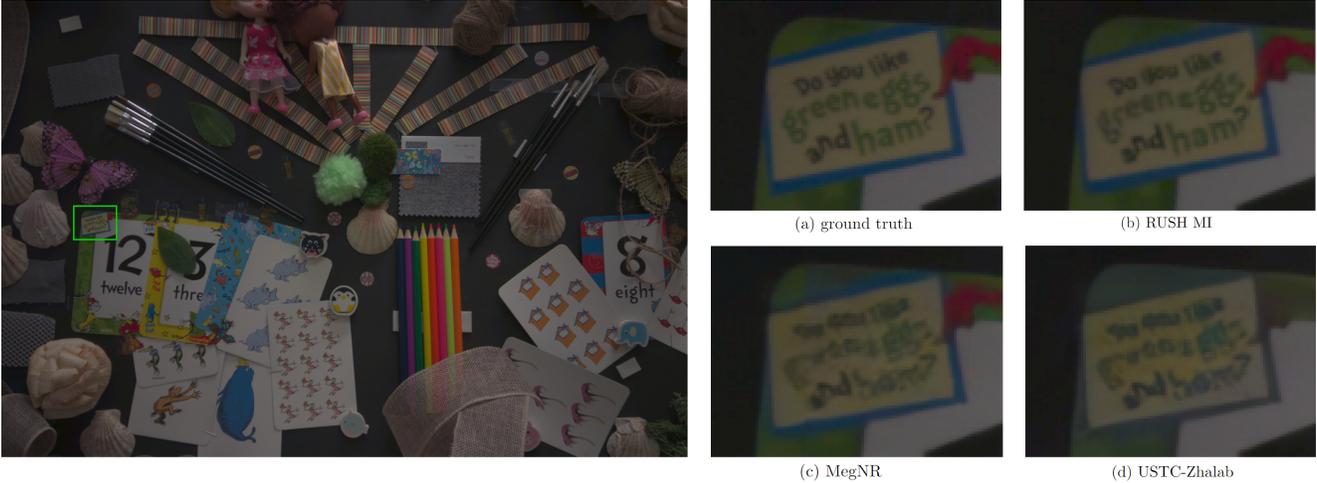
Figure 4. Qualitative image quality (IQ) comparison. The results of one of the test scenes (42dB) are shown. While the top three fusion methods achieve high objective IQ metrics in Table 1, texts on the card are slightly blurred in (b) and are barely interpretable in (c) and (d). The RGB images are obtained by using the ISP in Fig. 2, and its code is provided to participants.



Figure 5. Qualitative image quality (IQ) comparison. The results of one of the test scenes (42dB) are shown. Detail loss or false color in the top three methods in Table 1 can be found when compared with the ground truth. The mesh of the chair is over-smoothed to different extents in (c) and (d) and some false color can be found in (b). The RGB images are obtained by using the ISP in Fig. 2, and its code is provided to participants.

higher, power-efficient fusion algorithms are highly desirable.

## 4. Challenge Methods

This section describes the solutions submitted by all teams participating in the final stage of the MIPI 2023 RGBW Joint Fusion and Denoise Challenge.

### 4.1. RUSH MI

This team presents an end-to-end joint remosaic and denoise model, referred to as DEDD. As illustrated in Fig. 6, the DEDD model is composed of three components: a denoising model, a main network, and a differentiable ISP

model. Specifically, in the first part, they employed a basic UNet [10], which incorporated two downsampling operations. In the second part, they utilized the state-of-the-art model in the low-level domain, NAFNet [4]. The NAFNet contains the 4-level encoder-decoder and bottleneck. For the encoder, the numbers of NAFNet's blocks for each level are 2, 4, 8, and 24. For the decoder, the numbers of NAFNet's blocks for the 4 levels are all 2. In addition, the number of NAFNet's blocks for the bottleneck is 12. In the third part, they reformulated the BLC, WBC, GAMMA, and CCM modules in the conventional ISP pipeline into differentiable models, and adopted the officially provided demosaic model as the demosaic module. In the training phase, the clean model is used as a guidance for boosting

the noisy restoration performance. The images were randomly cut into 256×256 patches, with a batch size of 64. The optimizer used was AdamW [9], and the initial learning rate was set to 0.001, which was reduced by half every 5000 iterations. The training process is divided into two stages; initially, the denoising network is trained for 40k iterations, after which the parameters of the denoiser are fixed and the demosaic network is trained. The loss function utilized is the L1 loss. When training the demosaic network, two supervision signals are incorporated: the constraint of the Bayer domain and the constraint of the RGB domain. The model training is completed after retraining 80K iterations in an end-to-end training mode.



Figure 6. The model architecture of RUSH MI.

## 4.2. MegNR



Figure 7. The model architecture of MegNR.

This team proposes a novel two-stream pipeline based NAF-blocks [4] for RGBW joint fusion and denoise task, as shown in Fig. 7. Inspired by high signal-to-noise ratio white pixels, they introduce a new module called W-guided dynamic convolution(WGDC), which aggregates the spatial and channel attributes of the white-pixel feature, guides the dynamic change of network capability according to the white-pixel characteristics. Moreover, the authors design a method for synthesizing RGBW data, which effectively reduces the gap between synthetic data and real data. They use the official simple ISP code to transfer the standard Bayer to the RGB image for loss calculation, which consists of PSNR, SSIM, LPIPS [22]. For training, the authors randomly crop the training images into 128x128-sized patches with the 8-sized batches. Bayer Preserving Augmentation [8], Cutmix [19] are used for data augmentation.

They use cosine decay strategy to decrease the learning rate to $1 \times 10^{-7}$ with the initial learning rate $1 \times 10^{-3}$ and the entire training costs about 5 days and converges after $4 \times 10^{-6}$ iters. In the final inference stage, Test-time Augmentation [11] is used to get final result.
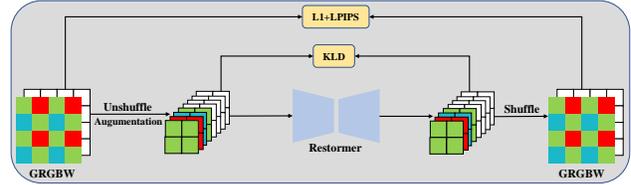
## 4.3. USTC-Zhalab



Figure 8. The model architecture of USTC-Zhalab.

This team proposes an RGBW fusion and denoising method based on the existing image restoration model Restormer [20], as shown in Fig. 8. During training, the Pixel-Unshuffle [13] is firstly applied to RGBW images to split them from 2 channels into 8 channels. Then, the 8-channel RGBW images are fed into the Restormer, obtaining the output of 8 channels. Finally, the Pixel-Shuffle [12] restores the output of 8 channels to the standard Bayer format. The training loss function consists of L1 loss, KLD loss, and LPIPS loss [22], calculated on the output of 8 channels and GT Bayer. Moreover, they also utilize three data augmentation strategies for training, i.e., Bayer Preserving Augmentation [8], Cutmix [19], and Mixup [21]. The whole network is trained for $3 \times 10^5$ iterations. The learning rate is decayed from $2 \times 10^{-4}$ to $1 \times 10^{-6}$ with a CosineAnnealing schedule. The training batch size and patch size are set to 8 and 224. The Self-ensemble strategy, Test-time Augmentation [11], and Test-time Local Converter [6] are applied during the testing phase. The testing batch size and patch size are set to 1 and 224.
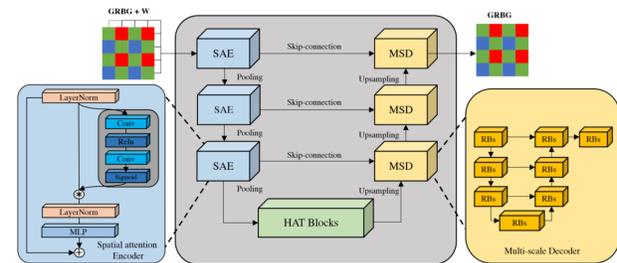
## 4.4. VIDAR



Figure 9. The model architecture of VIDAR.

This team proposes a multi-scale hybrid attention network for RGBW Fusion and denoising task as shown in Fig. 9. Inspired by Restormer [20] and HAT [5], the proposed method employs the Spatial Attention Module as the

decoder (SAE), which is decoded by the Multi-Scale Decoder (MSD) via skip-connections. The hybrid attention transformer (HAT) is also used in this strategy to consider both global and local information. The combination of these techniques enables efficient processing of high-resolution images as well as the extraction of significant information. The training process includes two stages. In the stage 1, the network is trained with the patches of size $256 \times 256$. The batch size is set to 12 and the optimizer is ADAM by setting of $\beta_1 = 0.9$ and $\beta_2 = 0.999$. The learning rate is initialized as $10^{-4}$ and it is decayed by a factor of $0.5$ at $50,000$, $80,000$, and $100,000$ iterations. Thanks to the network in stage 1 to greatly improve Bayer by using the white-channel information. The stage 2 finetunes the network from stage 1. In stage 2, the network is trained with the $L_1$ loss and LPIPS loss on RGB domain. The learning rate is initialized as $5 \times 10^{-5}$ and decayed by a factor of $0.5$ at $10,000$, $20,000$, and $40,000$ iterations. In the testing phase, the results are from the two stages to achieve the best performance. To be specific, the self-ensemble is used in the period of testing single model to get a better result. The input frame is flipped and regard it as another input. Then an inverse transform is applied to the corresponding output. An average of the transformed output and original output will be the self-ensemble result. The final result is a mixed results of the outputs from different iterations.

## 5. Conclusions

This report reviewed and summarized the methods and results of the RGBW Fusion challenge in the 2nd Mobile Intelligent Photography and Imaging workshop (MIPI 2023) held in conjunction with CVPR 2023. The participants were provided with a high-quality dataset for RGBW fusion and denoising. The four submissions leverage learning-based methods and achieve promising results. We are excited to see so many submissions within such a short period, and we look forward to more research in this area.

## 6. Acknowledgements

## References

[1] Camon 19 pro. https://www.tecno-mobile.com/phones/product-detail/product/camon-19-pro-5g. 1

[2] Oppo unveils multiple innovative imaging technologies. https://www.oppo.com/en/newsroom/press/oppo-future-imaging-technology-launch/. 1

[3] vivo x80 is the only vivo smartphone with a sony imx866 sensor: The world's first rgbw bottom sensors. https://www.vivoglobal.ph/vivo-X80-is-the-only-vivo-smartphone-with-a-Sony-IMX866-Sensor-The-Worlds-First-RGBW-Bottom-Sensors/. 1

[4] Liangyu Chen, Xiaojie Chu, Xiangyu Zhang, and Jian Sun. Simple baselines for image restoration. *arXiv preprint arXiv:2204.04676*, 2022. 4, 5

[5] Xiangyu Chen, Xintao Wang, Jiantao Zhou, Yu Qiao, and Chao Dong. Activating more pixels in image super-resolution transformer, 2023. 5

[6] Xiaojie Chu, Liangyu Chen, Chengpeng Chen, and Xin Lu. Revisiting global statistics aggregation for improving image restoration. *arXiv preprint arXiv:2112.04491*, 2021. 5

[7] Ruicheng Feng, Chongyi Li, Shangchen Zhou, Wenxiu Sun, Qingpeng Zhu, Jun Jiang, Qingyu Yang, Chen Change Loy, Jinwei Gu, Yurui Zhu, et al. Mipi 2022 challenge on under-display camera image restoration: Methods and results. In *In Proceedings of the European Conference on Computer Vision (ECCV) 2022 Workshops*, 2023. 2

[8] Jiaming Liu, Chi-Hao Wu, Yuzhi Wang, Qin Xu, Yuqian Zhou, Haibin Huang, Chuan Wang, Shaofan Cai, Yifan Ding, Haoqiang Fan, et al. Learning raw image denoising with bayer pattern unification and bayer preserving augmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 0–0, 2019. 5

[9] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 5

[10] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18*, pages 234–241. Springer, 2015. 4

[11] Divya Shanmugam, Davis Blalock, Guha Balakrishnan, and John Guttag. Better aggregation in test-time augmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1214–1223, 2021. 5

[12] Wenzhe Shi, Jose Caballero, Ferenc Huszár, Johannes Totz, Andrew P Aitken, Rob Bishop, Daniel Rueckert, and Zehan Wang. Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1874–1883, 2016. 5

[13] Bin Sun, Yulun Zhang, Songyao Jiang, and Yun Fu. Hybrid pixel-unshuffled network for lightweight image super-resolution. *arXiv preprint arXiv:2203.08921*, 2022. 5

[14] Wenxiu Sun, Qingpeng Zhu, Chongyi Li, Ruicheng Feng, Shangchen Zhou, Jun Jiang, Qingyu Yang, Chen Change Loy, Jinwei Gu, Dewang Hou, et al. Mipi 2022 challenge on rgb+ tof depth completion: Dataset and report. In *In Proceedings of the European Conference on Computer Vision (ECCV) 2022 Workshops*, 2023. 2

[15] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to

structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. 2, 3

[16] Qingyu Yang, Guang Yang, Jun Jiang, Chongyi Li, Ruicheng Feng, Shangchen Zhou, Wenxiu Sun, Qingpeng Zhu, Chen Change Loy, Jinwei Gu, et al. Mipi 2022 challenge on quad-bayer re-mosaic: Dataset and report. In *In Proceedings of the European Conference on Computer Vision (ECCV) 2022 Workshops*, 2023. 2

[17] Qingyu Yang, Guang Yang, Jun Jiang, Chongyi Li, Ruicheng Feng, Shangchen Zhou, Wenxiu Sun, Qingpeng Zhu, Chen Change Loy, Jinwei Gu, et al. Mipi 2022 challenge on rgbw sensor fusion: Dataset and report. In *In Proceedings of the European Conference on Computer Vision (ECCV) 2022 Workshops*, 2023. 2

[18] Qingyu Yang, Guang Yang, Jun Jiang, Chongyi Li, Ruicheng Feng, Shangchen Zhou, Wenxiu Sun, Qingpeng Zhu, Chen Change Loy, Jinwei Gu, et al. Mipi 2022 challenge on rgbw sensor re-mosaic: Dataset and report. In *In Proceedings of the European Conference on Computer Vision (ECCV) 2022 Workshops*, 2023. 2

[19] Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6023–6032, 2019. 5

[20] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, and Ming-Hsuan Yang. Restormer: Efficient transformer for high-resolution image restoration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5728–5739, 2022. 5

[21] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017. 5

[22] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 2, 3, 5

## A. Teams and Affiliations

### RUSH MI

**Title**:
Decoupled End-to-End Remosaic and Denoise Mode
**Members**:
Hongyuan Yu[1] (yuhongyuan@xiaomi.com),
Yuqinq Liu[2], Weichen Yu[3], Lin Ge[1], Xiaolin Zhang[1], Qi Jia[2], Heng Zhang[1], Xuanwu Yin[1], Kunlong Zuo[1]
**Affiliations**:
[1] Multimedia Department, Xiaomi Inc.,
[2] Dalian University of Technology,
[3] Institute of Automation, Chinese Academy of Sciences

### MegNR

**Title**:
WGDCNet: W-guided Dynamic Convolution Network for RGBW Fusion Image
**Members**:
Qi Wu[1] (wuqi02@megvii.com),
Wenjie Lin[1], Ting Jiang[1], Chengzhi Jiang[1], Mingyan Han[1], Xinpeng Li[1], Jinting Luo[1], Lei Yu[1], Haoqiang Fan[1], Shuaicheng Liu[2,1*]
**Affiliations**:
[1] Megvii Technology,
[2] University of Electronic Science and Technology of China (UESTC)

### USTC-Zhalab

**Title**:
Restormer-based RGBW Joint Fusion and Denoise
**Members**:
Kunyu Wang (kunyuwang@mail.ustc.edu.cn),
Chengzhi Cao
**Affiliations**:
University of Science and Technology of China

### VIDAR

**Title**:
Multi-Scale Hybrid Attention Network for RGBW Fusion and Denoising
**Members**:
Yuanshen Guan (guanys@mail.ustc.edu.cn),
Jiyuan Xia, Ruikang Xu, Mingde Yao, Zhiwei Xiong
**Affiliations**:
University of Science and Technology of China