

Dynamic Dual Trainable Bounds
for Ultra-low Precision Super-Resolution
Networks

Yunshan Zhong^{1,2}, Mingbao Lin³, Xunchao Li², Ke Li³, Yunhang Shen³,
Fei Chao^{1,2}, Yongjian Wu³, and Rongrong Ji^{1,2(B)}

¹Institute of Artificial Intelligence, Xiamen University, Xiamen, China
zhongyunshan@stu.xmu.edu.cn, {fchao, rrji}@xmu.edu.cn

²MAC Lab, School of Informatics, Xiamen University, Xiamen, China
lixunchao@stu.xmu.edu.cn

³Tencent Youtu Lab, Shanghai, China
littlekenwu@tencent.com

Abstract. Light-weight super-resolution (SR) models have received considerable attention for their serviceability in mobile devices. Many efforts employ network quantization to compress SR models. However, these methods suffer from severe performance degradation when quantizing the SR models to ultra-low precision (e.g., 2-bit and 3-bit) with the low-cost layer-wise quantizer. In this paper, we identify that the performance drop comes from the contradiction between the layer-wise symmetric quantizer and the highly asymmetric activation distribution in SR models. This discrepancy leads to either a waste on the quantization levels or detail loss in reconstructed images. Therefore, we propose a novel activation quantizer, referred to as Dynamic Dual Trainable Bounds (DDTB), to accommodate the asymmetry of the activations.

Specifically, DDTB innovates in: 1) A layer-wise quantizer with trainable upper and lower bounds to tackle the highly asymmetric activations. 2) A dynamic gate controller to adaptively adjust the upper and lower bounds at runtime to overcome the drastically varying activation ranges over different samples. To reduce the extra overhead, the dynamic gate controller is quantized to 2-bit and applied to only part of the SR networks according to the introduced dynamic intensity. Extensive experiments demonstrate that our DDTB exhibits significant performance improvements in ultra-low precision. For example, our DDTB achieves a 0.70 dB PSNR increase on Urban100 benchmark when quantizing EDSR to 2-bit and scaling up output images to $\times 4$. Code is at <https://github.com/zysxmu/DDTB>.

Keywords: Super-resolution
· Network quantization · Dual trainable
bounds · Dynamic gate controller

Supplementary Information The online version contains supplementary material available at https://doi.org/10.1007/978-3-031-19797-0_1.

© The Author(s), under exclusive license to Springer Nature Switzerland AG 2022
S. Avidan et al. (Eds.): ECCV 2022, LNCS 13678, pp. 1–18, 2022. <https://doi.org/10.1007/978-3-031-19797-0>

OSFormer: One-Stage Camouflaged

Instance Segmentation with Transformers

Jialun Pei¹, Tianyang Cheng², Deng-Ping Fan^{3(B)}, Hengshuang Zhao²,
Chuanbo Chen², and Luc Van Gool³

¹School of Computer Science and Technology, HUST, Wuhan, China

²School of Software Engineering, HUST, Wuhan, China

³Computer Vision Lab, ETH Zurich, Zurich, Switzerland
dengpfan@gmail.com

Abstract. We present OSFormer, the first one-stage transformer framework for camouflaged instance segmentation (CIS). OSFormer is based on two key designs. First, we design a location-sensing transformer (LST) to obtain the location label and instance-aware parameters by introducing the location-guided queries and the blend-convolution feed-forward network. Second, we develop a coarse-to-fine fusion (CFF) to merge diverse context information from the LST encoder and

CNN backbone. Coupling these two components enables OSFormer to efficiently blend local features and long-range context dependencies for predicting camouflaged instances. Compared with two-stage frameworks, our OSFormer reaches 41% AP and achieves good convergence efficiency without requiring enormous training data, i.e., only 3,040 samples under 60 epochs. Code link: <https://github.com/PJLallen/OSFormer>.

Keywords: Camouflage

• Instance segmentation • Transformer

1

Highly Accurate Dichotomous Image Segmentation

Xuebin Qin¹, Hang Dai¹, Xiaobin Hu², Deng-Ping Fan^{3(B)},
Ling Shao⁴, and Luc Van Gool³

¹MBZUAI, Abu Dhabi, UAE

xuebin@ualberta.ca, hang.dai@mbzuai.ac.ae

²Tencent Youtu Lab, Shanghai, China

xiaobin.hu@tum.de

³ETH Zurich, Zurich, Switzerland

dengpfan@gmail.com, vangool@vision.ee.ethz.ch

⁴Terminus Group, Beijing, China

ling.shao@ieee.org

Abstract. We present a systematic study on a new task called dichotomous image segmentation (DIS), which aims to segment highly accurate objects from natural images. To this end, we collected the first large-scale DIS dataset, called DIS5K, which contains 5,470 high-resolution (e.g., 2K, 4K or larger) images covering camouflaged, salient, or meticulous objects in various backgrounds. DIS is annotated with extremely fine-grained labels. Besides, we introduce a simple intermediate supervision baseline (IS-Net) using both feature-level and mask-level guidance for DIS model training. IS-Net outperforms various cutting-edge baselines on the proposed DIS5K, making it a general self-learned supervision network that can facilitate future research in DIS. Further, we design a new metric called human correction efforts (HCE) which approximates the number of mouse clicking operations required to correct the false positives and false negatives. HCE is utilized to measure the gap between models and real-world applications and thus can complement existing metrics. Finally, we conduct the largest-scale benchmark, evaluating 16 representative segmentation models, providing a more insightful discussion regarding object complexities, and showing several potential applications (e.g., background removal, art design, 3D reconstruction). Hoping these efforts can open up promising directions for both academic and industries. Project page: <https://xuebinqin.github.io/dis/index.html>.

Keywords: Dichotomous image segmentation

• High resolution •

Metric

We would like to thank Jiayi Zhu for his efforts in re-organizing the dataset and codes.

Supplementary Information The online version contains supplementary material available at https://doi.org/10.1007/978-3-031-19797-0_3.

© The Author(s), under exclusive license to Springer Nature Switzerland AG 2022
S. Avidan et al. (Eds.): ECCV 2022, LNCS 13678, pp. 38–56, 2022. <https://doi.org/10.1007/978-3-031-19797-0>

–

Boosting Supervised Dehazing Methods

via Bi-level Patch Reweighting

Xingyu Jiang¹, Hongkun Dou¹, Chengwei Ful, Bingquan Dai¹, Tianrun Xu²,
and Yue Deng^{1(B)}

¹School of Astronautics, Beihang University, Beijing, China

ydeng@buaa.edu.cn

2North China University of Technology, Beijing, China

Abstract. Natural images can suffer from non-uniform haze distributions in different regions. However, this important fact is hardly considered in existing supervised dehazing methods, in which all training patches are accounted for equally in the loss design. These supervised methods may fail in making promising recoveries on some regions contaminated by heavy hazes. Therefore, for a more reasonable dehazing losses design, the varying importance of different training patches should be taken into account. Such rationale is exactly in line with the process of human learning that difficult concepts always require more practice in learning. To this end, we propose a bi-level dehazing (BILD) framework by designing an internal loop for weighted supervised dehazing and an external loop for training patch reweighting. With simple derivations, we show the gradients of BILD exhibit natural connections with policy gradient and can thus explain the BILD objective by the rewarding mechanism in reinforcement learning. The BILD is not a new dehazing method per se, it is better recognized as a flexible framework that can seamlessly work with general supervised dehazing approaches for their performance boosting.

Keywords: Single image dehazing

·Bi-level optimization ·Visual

importance ·Deep learning

1

Flow-Guided Transformer for Video

Inpainting

Kaidong Zhang¹, Jingjing Fu^{2(B)}, and Dong Liu¹

¹University of Science and Technology of China, Hefei, China

richu@mail.ustc.edu.cn ,dongeliu@ustc.edu.cn

²Microsoft Research Asia, Beijing, China

jifu@microsoft.com

Abstract. We propose a flow-guided transformer, which innovatively leverage the motion discrepancy exposed by optical flows to instruct the attention retrieval in transformer for high fidelity video inpainting. More specially, we design a novel flow completion network to complete the corrupted flows by exploiting the relevant flow features in a local temporal window. With the completed flows, we propagate the content across video frames, and adopt the flow-guided transformer to synthesize the rest corrupted regions. We decouple transformers along temporal and spatial dimension, so that we can easily integrate the locally relevant completed flows to instruct spatial attention only. Furthermore, we design a flow-reweight module to precisely control the impact of completed flows on each spatial transformer. For the sake of efficiency, we introduce window partition strategy to both spatial and temporal transformers. Especially in spatial transformer, we design a dual perspective spatial MHSA, which integrates the global tokens to the window-based attention. Extensive experiments demonstrate the effectiveness of the proposed method qualitatively and quantitatively. Codes are available at <https://github.com/hitachinsk/FGT>.

Keywords: Video inpainting

·Optical flow ·Transformer

1

Shift-Tolerant Perceptual Similarity

Metric

Abhijay Ghildyal^(B) and Feng Liu

Portland State University, Portland, OR 97201, USA

{abhijay, fliu}@pdx.edu

Abstract. Existing perceptual similarity metrics assume an image and its reference are well aligned. As a result, these metrics are often sensitive to

a small alignment error that is imperceptible to the human eyes. This paper studies the effect of small misalignment, specifically a small shift between the input and reference image, on existing metrics, and accordingly develops a shift-tolerant similarity metric. This paper builds upon LPIPS, a widely used learned perceptual similarity metric, and explores architectural design considerations to make it robust against imperceptible misalignment. Specifically, we study a wide spectrum of neural network elements, such as anti-aliasing filtering, pooling, striding, padding, and skip connection, and discuss their roles in making a robust metric. Based on our studies, we develop a new deep neural network-based perceptual similarity metric. Our experiments show that our metric is tolerant to imperceptible shifts while being consistent with the human similarity judgment. Code is available at <https://tinyurl.com/5n85r28r>.

Keywords: Perceptual similarity metric

•Image quality assessment

1

Perception-Distortion Balanced ADMM

Optimization for Single-Image

Super-Resolution

Yuehan Zhang¹, Bo Ji¹, Jiahua Yao², and Angela Yao^{1(B)}

¹National University of Singapore, Singapore, Singapore

{zyuehan,jibo,ayao}@comp.nus.edu.sg

²HiSilicon Technologies, Shanghai, China

hao.jia@huawei.com

Abstract. In image super-resolution, both pixel-wise accuracy and perceptual fidelity are desirable. However, most deep learning methods only achieve high performance in one aspect due to the perception-distortion trade-off, and works that successfully balance the trade-off rely on fusing results from separately trained models with ad-hoc post-processing. In this paper, we propose a novel super-resolution model with a low-frequency constraint (LFC-SR), which balances the objective and perceptual quality through a single model and yields super-resolved images with high PSNR and perceptual scores. We further introduce an ADMM-based alternating optimization method for the non-trivial learning of the constrained model. Experiments showed that our method, without cumbersome post-processing procedures, achieved the state-of-the-art performance. The code is available at <https://github.com/Yuehan717/PDASR>.

Keywords: Image super-resolution

•Perception-distortion trade-off •

Constrained optimization

1

VQFR: Blind Face Restoration

with Vector-Quantized Dictionary

and Parallel Decoder

Yuchao Gu^{1,2}, Xintao Wang², Liangbin Xie^{2,5}, Chao Dong^{4,5}, Genliang Ying², and Ming-Ming Cheng^{1(B)}

¹TMCC, CS, Nankai University, Tianjin, China

cmm@nankai.edu.cn

²ARC Lab, Tencent PCG, Beijing, China

³Platform Technologies, Tencent Online Video, Beijing, China

⁴Shanghai AI Laboratory, Beijing, China

⁵Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences, Beijing, China

<https://github.com/TencentARC/VQFR/>

Abstract. Although generative facial prior and geometric prior have recently demonstrated high-quality results for blind face restoration, producing fine-grained facial details faithful to inputs remains a challenging problem. Motivated by the classical dictionary-based methods and the

recent vector quantization (VQ) technique, we propose a VQ-based facerestoration method - VQFR. VQFR takes advantage of high-quality low-level feature banks extracted from high-quality faces and can thus help recover realistic facial details. However, the simple application of the VQ codebook cannot achieve good results with faithful details and identity preservation. Therefore, we further introduce two special network designs. 1). We first investigate the compression patch size in the VQ codebook and find that the VQ codebook designed with a proper compression patch size is crucial to balance the quality and fidelity. 2). To further fuse low-level features from inputs while not "contaminating" the realistic details generated from the VQ codebook, we proposed a parallel decoder consisting of a texture decoder and a main decoder. Those two decoders then interact with a texture warping module with deformable convolution. Equipped with the VQ codebook as a facial detail dictionary and the parallel decoder design, the proposed VQFR can largely enhance the restored quality of facial details while keeping the fidelity to previous methods.

Keywords: Blind face restoration

• Vector quantization • Parallel decoder

X. Wang—Project lead. Y. Gu—is an intern in ARC Lab, Tencent PCG.

Supplementary Information The online version contains supplementary material available at https://doi.org/10.1007/978-3-031-19797-0_8.

© The Author(s), under exclusive license to Springer Nature Switzerland AG 2022
S. Avidan et al. (Eds.): ECCV 2022, LNCS 13678, pp. 126–143, 2022. <https://doi.org/10.1007/978-3-031-19797-0>

Uncertainty Learning in Kernel
Estimation for Multi-stage Blind Image
Super-Resolution

Zhenxuan Fang¹, Weisheng Dong^{1(B)}, Xin Li², Jinjian Wu¹, Leida Li¹,
and Guangming Shi¹

¹School of Artificial Intelligence, Xidian University, Xi'an, China
zxfang@stu.xidian.edu.cn, {wsdong, jinjian.wu}@mail.xidian.edu.cn,
{ldli, gmshi}@xidian.edu.cn

²Lane Department of CSEE, West Virginia University, Morgantown, WV, USA
xin.li@mail.wvu.edu

Abstract. Conventional wisdom in blind super-resolution (SR) first estimates the unknown degradation from the low-resolution image and then exploits the degradation information for image reconstruction. Such sequential approaches suffer from two fundamental weaknesses - i.e., the lack of robustness (the performance drops when the estimated degradation is inaccurate) and the lack of transparency (network architectures are heuristic without incorporating domain knowledge). To address these issues, we propose a joint Maximum a Posteriori (MAP) approach for estimating the unknown kernel and high-resolution image simultaneously.

Our method first introduces uncertainty learning in the latent space when estimating the blur kernel, aiming at improving the robustness to the estimation error. Then we propose a novel SR network by unfolding the joint MAP estimator with a learned Laplacian Scale Mixture (LSM) prior and the estimated kernel. We have also developed a novel approach of estimating both the scale prior coefficient and the local means of the LSM model through a deep convolutional neural network (DCNN). All parameters of the MAP estimation algorithm and the DCNN parameters are jointly optimized through end-to-end training. Extensive experiments on both synthetic and real-world images show that our method achieves state-of-the-art performance for the task of blind image SR.

1

Learning Spatio-Temporal Downsampling for Effective Video Upscaling

Xiaoyu Xiang¹(B), Yapeng Tian², Vijay Rengarajan¹,
Lucas D. Young¹, Bo Zhang¹, and Rakesh Ranjan¹
¹Meta Reality Labs, Menlo Park, USA
{xiangxiaoyu,apvijay,bozhufri,rakeshr}@fb.com, lucasyoung482@gmail.com
²University of Texas at Dallas, Richardson, USA
tianyapeng92@gmail.com

Abstract. Downsampling is one of the most basic image processing operations. Improper spatio-temporal downsampling applied on videos can cause aliasing issues such as moiré patterns in space and the wagon-wheel effect in time. Consequently, the inverse task of upscaling a low-resolution, low frame-rate video in space and time becomes a challenging ill-posed problem due to information loss and aliasing artifacts. In this paper, we aim to solve the space-time aliasing problem by learning a spatio-temporal downsampler. Towards this goal, we propose a neural network framework that jointly learns spatio-temporal downsampling and upsampling. It enables the downsampler to retain the key patterns of the original video and maximizes the reconstruction performance of the upsampler. To make the downsampling results compatible with popular image and video storage formats, the downsampling results are encoded to uint8 with a differentiable quantization layer. To fully utilize the space-time correspondences, we propose two novel modules for explicit temporal propagation and space-time feature rearrangement. Experimental results show that our proposed method significantly boosts the space-time reconstruction quality by preserving spatial textures and motion patterns in both downsampling and upscaling. Moreover, our framework enables a variety of applications, including arbitrary video resampling, blurry frame reconstruction, and efficient video storage.

Keywords: Downsampling

·Anti-aliasing ·Video upscaling

1

Learning Local Implicit Fourier Representation for Image Warping

Jaewon Lee¹, Kwang Pyo Choi², and Kyong Hwan Jin¹(B)
¹Daegu Gyeongbuk Institute of Science and Technology (DGIST), Daegu, Korea
{ljw3136,kyong.jin}@dgist.ac.kr
²Samsung Electronics, Suwon-si, Korea
kp5.choi@samsung.com

Abstract. Image warping aims to reshape images defined on rectangular grids into arbitrary shapes. Recently, implicit neural functions have shown remarkable performances in representing images in a continuous manner. However, a standalone multi-layer perceptron suffers from learning high-frequency Fourier coefficients. In this paper, we propose a local texture estimator for image warping (LTEW) followed by an implicit neural representation to deform images into continuous shapes. Local textures estimated from a deep super-resolution (SR) backbone are multiplied by locally-varying Jacobian matrices of a coordinate transformation to predict Fourier responses of a warped image. Our LTEW-based neural function outperforms existing warping methods for asymmetric-scale SR and homography transform. Furthermore, our algorithm well generalizes arbitrary coordinate transformations, such as homography transform with a large magnification factor and equirectangular projection (ERP) perspective transform, which are not provided in training. Our source code is available at <https://github.com/jaewon-lee-b/ltew>.

Keywords: Image warping

·Implicit neural representation ·Fourier
features ·Jacobian ·Homography transform ·Equirectangular
projection (ERP)

SepLUT: Separable Image-Adaptive Lookup Tables for Real-Time Image Enhancement

Canqian Yang¹, Meiguang Jin², Yixun (B), Rui Zhang¹, Ying Chen²,
and Huaida Liu²

¹MoE Key Lab of Artificial Intelligence, AI Institute,
Shanghai Jiao Tong University, Shanghai, China

{charles.young,xuyi,zhang rui}@sjtu.edu.cn

²Alibaba Group, Hangzhou, China

{meiguang.jmg,yingchen,liuhuaيدا.1hd}@alibaba-inc.com

Abstract. Image-adaptive lookup tables (LUTs) have achieved great success in real-time image enhancement tasks due to their high efficiency for modeling color transforms. However, they embed the complete transform, including the color component-independent and the component-correlated parts, into only a single type of LUTs, either 1D or 3D, in a coupled manner. This scheme raises a dilemma of improving model expressiveness or efficiency due to two factors. On the one hand, the 1D LUTs provide high computational efficiency but lack the critical capability of color components interaction. On the other, the 3D LUTs present enhanced component-correlated transform capability but suffer from heavy memory footprint, high training difficulty, and limited cell utilization. Inspired by the conventional divide-and-conquer practice in the image signal processor, we present SepLUT (separable image-adaptive lookup table) to tackle the above limitations. Specifically, we separate a single color transform into a cascade of component-independent and component-correlated sub-transforms instantiated as 1D and 3D LUTs, respectively. In this way, the capabilities of two sub-transforms can facilitate each other, where the 3D LUT complements the ability to mix up color components, and the 1D LUT redistributes the input colors to increase the cell utilization of the 3D LUT and thus enable the use of a more lightweight 3D LUT. Experiments demonstrate that the proposed method presents enhanced performance on photo retouching benchmark datasets than the current state-of-the-art and achieves real-time processing on both GPUs and CPUs.

C. Yang and M. Jin—Equal contribution.

Work partially done during an internship of C. Yang at Alibaba Group.

Supplementary Information The online version contains supplementary material available at https://doi.org/10.1007/978-3-031-19797-0_12.

© The Author(s), under exclusive license to Springer Nature Switzerland AG 2022
S. Avidan et al. (Eds.): ECCV 2022, LNCS 13678, pp. 201–217, 2022. https://doi.org/10.1007/978-3-031-19797-0_12

_1

Blind Image Decomposition

Junlin Han^{1,2(B)}, Weihao Lil, Pengfei Fan^{gl,2}, Chunyi Sun², Jie Hong^{gl,2},
2,

Mohammad Ali Armin¹, Lars Petersson¹, and Hongdong Li²

¹Data61-CSIRO, Sydney, Australia

junlin.han@data61.csiro.au

²Australian National University, Canberra, Australia

Abstract. We propose and study a novel task named Blind Image Decomposition (BID), which requires separating a superimposed image into constituent underlying images in a blind setting, that is, both the source components involved in mixing as well as the mixing mechanism are unknown. For example, rain may consist of multiple components, such as rain streaks, raindrops, snow, and haze. Rainy images can be treated as an arbitrary combination of these components, some of them or all of them. How to decompose superimposed images, like rainy images,

into distinct source components is a crucial step toward real-world vision systems. To facilitate research on this new task, we construct multiple benchmark datasets, including mixed image decomposition across multiple domains, real-scenario deraining, and joint shadow/reflection/watermark removal. Moreover, we propose a simple yet general Blind Image Decomposition Network (BIDeN) to serve as a strong baseline for future work. Experimental results demonstrate the tenability of our benchmarks and the effectiveness of BIDeN.

Codes and datasets are available at GitHub .

Keywords: Image decomposition

·Low-level vision ·Rain removal

1

MuLUT: Cooperating Multiple Look-Up

Tables for Efficient Image

Super-Resolution

Jiacheng Li¹, Chang Chen², Zhen Cheng¹, and Zhiwei Xiong^{1(B)}

¹University of Science and Technology of China, Hefei, China

{jcleee,mywander}@mail.ustc.edu.cn, zwxiong@ustc.edu.cn

²Huawei Noah's Ark Lab, Beijing, China

chenchang25@huawei.com

Abstract. The high-resolution screen of edge devices stimulates a strong demand for efficient image super-resolution (SR). An emerging research, SR-LUT, responds to this demand by marrying the look-up table (LUT) with learning-based SR methods. However, the size of a single LUT grows exponentially with the increase of its indexing capacity. Consequently, the receptive field of a single LUT is restricted, resulting in inferior performance. To address this issue, we extend SR-LUT by enabling the cooperation of Multiple LUTs, termed MuLUT. Firstly, we devise two novel complementary indexing patterns and construct multiple LUTs in parallel. Secondly, we propose a re-indexing mechanism to enable the hierarchical indexing between multiple LUTs. In these two ways, the total size of MuLUT is linear to its indexing capacity, yielding a practical method to obtain superior performance. We examine the advantage of MuLUT on five SR benchmarks. MuLUT achieves a significant improvement over SR-LUT, up to 1.1 dB PSNR, while preserving its efficiency. Moreover, we extend MuLUT to address demosaicing of Bayer-patterned images, surpassing SR-LUT on two benchmarks by a large margin.

Keywords: Image super-resolution

·Look-up table ·Image

demosaicing

1

Learning Spatiotemporal

Frequency-Transformer for Compressed

Video Super-Resolution

Zhongwei Qiu^{1,2(B)}, Huan Yang³, Jianlong Fu³, and Dongmei Fu^{1,2}

¹University of Science and Technology Beijing, Beijing, China

qiuzhongwei@xs.ustb.edu.cn, fdm_ustb@ustb.edu.cn

²Shunde Graduate School of University of Science and Technology Beijing, Beijing, China

³Microsoft Research, Beijing, China

{huayan,jianf}@microsoft.com

Abstract. Compressed video super-resolution (VSR) aims to restore high-resolution frames from compressed low-resolution counterparts. Most recent VSR approaches often enhance an input frame by "borrowing" relevant textures from neighboring video frames. Although some progress has been made, there are grand challenges to effectively extract and transfer high-quality textures from compressed videos where most frames are u

usually highly degraded. In this paper, we propose a novel Frequency-Transformer for compressed video super-resolution (FTVSR) that conducts self-attention over a joint space-time-frequency domain. First, we divide a video frame into patches, and transform each patch into DCT spectral maps in which each channel represents a frequency band. Such a design enables a fine-grained level self-attention on each frequency band, so that real visual texture can be distinguished from artifacts, and further utilized for video frame restoration. Second, we study different self-attention schemes, and discover that a "divided attention" which conducts a joint space-frequency attention before applying temporal attention on each frequency band, leads to the best video enhancement quality. Experimental results on two widely-used video super-resolution benchmarks show that FTVSR outperforms state-of-the-art approaches on both uncompressed and compressed videos with clear visual margins. Code are available at <https://github.com/researchmm/FTVSR>.

Keywords: VSR

·Transformer ·Frequency learning ·Compression

1

Spatial-Frequency Domain Information

Integration for Pan-Sharpening

Man Zhou^{1,2}, Jie Huang¹, Keyu Yan^{1,2}, Hu Yu¹, Xueyang Fu¹, Aiping Liu¹, Xian Wei³, and Feng Zhao^{1(B)}

¹University of Science and Technology of China, Hefei, China

{manman,hj0117}@mail.ustc.edu.cn, fzhao956@ustc.edu.cn

²Hefei Institute of Physical Science, Chinese Academy of Sciences, Hefei, China

³MoE Engineering Research Center of Hardware/Software Co-design Technology and Application, East China Normal University, Shanghai, China

Abstract. Pan-sharpening aims to generate high-resolution multi-spectral (MS) images by fusing PAN images and low-resolution MS images. Despite its great advances, most existing pan-sharpening methods only work in the spatial domain and rarely explore the potential solutions in the frequency domain. In this paper, we first attempt to address pan-sharpening in both spatial and frequency domains and propose a Spatial-Frequency Information Integration Network, dubbed as SFIIN. To implement SFIIN, we devise a core building module tailored with pan-sharpening, consisting of three key components: spatial-domain information branch, frequency-domain information branch, and dual domain interaction. To be specific, the first employs the standard convolution to integrate the local information of two modalities of PAN and MS images in the spatial domain, while the second adopts deep Fourier transformation to achieve the image-wide receptive field for exploring global contextual information. Followed by, the third is responsible for facilitating the information flow and learning the complementary representation. We conduct extensive experiments to validate the effectiveness of the proposed network and demonstrate the favorable performance against other state-of-the-art methods.

Keywords: Pan-sharpening

·Spatial-frequency domain

1

Adaptive Patch Exiting for Scalable

Single Image Super-Resolution

Shizun Wang¹, Jiaming Liu^{2,4}, Kaixin Chen¹, Xiaoqi Li^{2,4}, Ming Lu^{3(B)}, and Yandong Guo⁴

¹Beijing University of Posts and Telecommunications, Beijing, China

wangshizun@bupt.edu.cn

²Peking University, Beijing, China

³Intel Labs China, Beijing, China

lu199192@gmail.com

4OPPO Research Institute, Shanghai, China

Abstract. Since the future of computing is heterogeneous, scalability is a crucial problem for single image super-resolution. Recent works try to train one network, which can be deployed on platforms with different capacities. However, they rely on the pixel-wise sparse convolution, which is not hardware-friendly and achieves limited practical speedup. As image can be divided into patches, which have various restoration difficulties, we present a scalable method based on Adaptive Patch Exiting (APE) to achieve more practical speedup. Specifically, we propose to train a regressor to predict the incremental capacity of each layer for the patch. Once the incremental capacity is below the threshold, the patch can exit at the specific layer. Our method can easily adjust the trade-off between performance and efficiency by changing the threshold of incremental capacity. Furthermore, we propose a novel strategy to enable the network training of our method. We conduct extensive experiments across various backbones, datasets and scaling factors to demonstrate the advantages of our method. Code is available at <https://github.com/littlepure2333/APE>.

Keywords: Single image super-resolution

· Scalability · Efficiency

1

Efficient Meta-Tuning for Content-Aware

Neural Video Delivery

Xiaoqi Li¹, Jiaming Liu^{1,2}, Shizun Wang³, Cheng Lyu³,
Ming Lu^{4(B)}, Yurong Chen⁴, Anban GYao⁴, Yandong Guo²,
and Shanghang Zhang^{1(B)}

¹Peking University, Beijing, China

shzhang.pku@gmail.com

²OPPO Research Institute, Beijing, China

³Beijing University of Posts and Telecommunications, Beijing, China

⁴Intel Labs China, Beijing, China

Abstract. Recently, Deep Neural Networks (DNNs) are utilized to reduce the bandwidth and improve the quality of Internet video delivery. Existing methods train corresponding content-aware super-resolution (SR) model for each video chunk on the server, and stream low-resolution (LR) video chunks along with SR models to the client. Although they achieve promising results, the huge computational cost of network training limits their practical applications. In this paper, we present a method named Efficient Meta-Tuning (EMT) to reduce the computational cost. Instead of training from scratch, EMT adapts a meta-learned model to the first chunk of the input video. As for the following chunks, it fine-tunes the partial parameters selected by gradient masking of previous adapted model. In order to achieve further speedup for EMT, we propose a novel sampling strategy to extract the most challenging patches from video frames. The proposed strategy is highly efficient and brings negligible additional cost. Our method significantly reduces the computational cost and achieves even better performance, paving the way for applying neural video delivery techniques to practical applications. We conduct extensive experiments based on various efficient SR architectures, including ESPCN, SRCNN, FSRNN and EDSR-1, demonstrating the generalization ability of our work. The code is released at <https://github.com/Neural-video-delivery/EMT-Pytorch-ECCV2022>.

Keywords: Neural video delivery

· Super-resolution · Meta learning

X. Li, J. Liu and S. Wang—Equal contribution.

Supplementary Information The online version contains supplementary material available at https://doi.org/10.1007/978-3-031-19797-0_18.

© The Author(s), under exclusive license to Springer Nature Switzerland AG 2022
S. Avidan et al. (Eds.): ECCV 2022, LNCS 13678, pp. 308–324, 2022. <https://doi.org/10.1007/978-3-031-19797-0>

Reference-Based Image Super-Resolution with Deformable Attention Transformer

Jie Zhang¹, Jingyun Liang¹, Kai Zhang¹, Yawei Li¹, Yulun Zhang^{1(B)},
Wenguan Wang¹, and Luc Van Gool^{1,2}

¹Computer Vision Lab, ETH Zürich, Zürich, Switzerland

{jie.zhang.cao, jingyun.liang, kai.zhang, yawei.li, yulun.zhang,
wenguan.wang, vangool }@vision.ee.ethz.ch

²KU Leuven, Leuven, Belgium

<https://github.com/caojiezhong/DATSR>

Abstract. Reference-based image super-resolution (RefSR) aims to exploit auxiliary reference (Ref) images to super-resolve low-resolution (LR) images. Recently, RefSR has been attracting great attention as it provides an alternative way to surpass single image SR. However, addressing the RefSR problem has two critical challenges: (i) It is difficult to match the correspondence between LR and Ref images when they are significantly different; (ii) How to transfer the relevant texture from Ref images to compensate the details for LR images is very challenging. To address these issues of RefSR, this paper proposes a deformable attention transformer, namely DATSR, with multiple scales, each of which consists of a texture feature encoder (TFE) module, a reference-based deformable attention (RDA) module and a residual feature aggregation (RFA) module. Specifically, TFE first extracts image transformation (e.g., brightness) insensitive features for LR and Ref images, RDA then can exploit multiple relevant textures to compensate more information for LR features, and RFA lastly aggregates LR features and relevant textures to get a more visually pleasant result. Extensive experiments demonstrate that our DATSR achieves state-of-the-art performance on benchmark datasets quantitatively and qualitatively.

Keywords: Reference-based image super-resolution

·Correspondence

matching ·Texture transfer ·Deformable attention transformer

1

Local Color Distributions Prior for Image Enhancement

Haoyuan Wang^(B), Ke Xu, and Rynson W.H. Lau

Department of Computer Science, City University of Hong Kong, Hong Kong,
People's Republic of China

hywang26-c@my.city.edu.hk

Abstract. Existing image enhancement methods are typically designed to address either the over- or under-exposure problem in the input image. When the illumination of the input image contains both over- and under-exposure problems, these existing methods may not work well. We observe from the image statistics that the local color distributions (LCDs) of an image suffering from both problems tend to vary across different regions of the image, depending on the local illuminations. Based on this observation, we propose in this paper to exploit these LCDs as a prior for locating and enhancing the two types of regions (i.e., over-/under-exposed regions). First, we leverage the LCDs to represent these regions, and propose a novel local color distribution embedded (LCDE) module to formulate LCDs in multi-scales to model the correlations across different regions. Second, we propose a dual-illumination learning mechanism to enhance the two types of regions. Third, we construct a new dataset to facilitate the learning process, by following the camera image signal processing (ISP) pipeline to render standard RGB images with both under-/over-exposures from raw data. Extensive experiments demonstrate that the proposed method outperforms existing state-of-the-art methods quantitatively and qualitatively. Codes and dataset are in <https://github.com/haoyuanwang/LCDE>

/hywang99.github.io/lcdpnet/ .

1

L-CoDer: Language-Based Colorization
with Color-Object Decoupling
Transformer

Zheng Chang¹, Shuchen Weng², Yulu³, Silu¹(B), and Boxin Shi²

¹School of Artificial Intelligence, Beijing University of Posts and
Telecommunications, Beijing, China

{zhengchang98, lisi}@bupt.edu.cn

²NERCVT, School of Computer Science, Peking University, Beijing, China

{shuchenweng, shiboxin}@pku.edu.cn

³International Digital Economy Academy, Shenzhen, China

liyu@idea.edu.cn

Abstract. Language-based colorization requires the colorized image to be consistent with the user-provided language caption. A most recent work proposes to decouple the language into color and object conditions in solving the problem. Though decent progress has been made, its performance is limited by three key issues. (i) The large gap between vision and language modalities using independent feature extractors makes it difficult to fully understand the language. (ii) The inaccurate language features are never refined by the image features such that the language may fail to colorize the image precisely. (iii) The local region does not perceive the whole image, producing global inconsistent colors. In this work, we introduce transformer into language-based colorization to tackle the aforementioned issues while keeping the language decoupling property.

Our method unifies the modalities of image and language, and further performs color conditions evolving with image features in a coarse-to-fine manner.

In addition, thanks to the global receptive field, our method is robust to the strong local variation. Extensive experiments demonstrate our method is able to produce realistic colorization and outperforms prior arts in terms of consistency with the caption.

1

From Face to Natural Image: Learning
Real Degradation for Blind Image
Super-Resolution

Xiaoming Li^{1,5}, Chaofeng Chen², Xianhui Lin³, Wangmeng Zuo^{1,4}(B),
and Lei Zhang⁵

¹Faculty of Computing, Harbin Institute of Technology, Harbin, China
wmzuo@hit.edu.cn

²S-Lab, Nanyang Technological University, Singapore, Singapore

³DAMO Academy, Alibaba Group, Shenzhen, China

⁴Peng Cheng Lab, Shenzhen, China

⁵Department of Computing, The Hong Kong Polytechnic University,
Hung Hom, Hong Kong

cslzhang@comp.polyu.edu.hk

Abstract. How to design proper training pairs is critical for super-resolving real-world low-quality (LQ) images, which suffers from the difficulties in either acquiring paired ground-truth high-quality (HQ) images or synthesizing photo-realistic degraded LQ observations. Recent works mainly focus on modeling the degradation with handcrafted or estimated degradation parameters, which are however incapable to model complicated real-world degradation types, resulting in limited quality improvement. Notably, LQ face images, which may have the same degradation process as natural images, can be robustly restored with photo-realistic textures by exploiting their strong structural priors. This motivates us to use the real-world LQ face images and their restored HQ counterparts to model the complex real-world degradation (namely ReDegNet), and then transfer it to HQ natural images to synthesize their realistic LQ

counterparts. By taking these paired HQ-LQ face images as inputs to explicitly predict the degradation-aware and content-independent representations, we could control the degraded image generation, and subsequently transfer these degradation representations from face to natural images to synthesize the degraded LQ natural images. Experiments show that our ReDegNet can well learn the real degradation process from face images. The restoration network trained with our synthetic pairs performs favorably against SOTAs. More importantly, our method provides a new way to handle the real-world complex scenarios by learning their degradation representations from the facial portions, which can be used to significantly improve the quality of non-facial areas. The source code is available at <https://github.com/csxmli2016/ReDegNet>.

Keywords: Real world degradation

•Blind image super-resolution

Supplementary Information The online version contains supplementary material available at https://doi.org/10.1007/978-3-031-19797-0_22.

©The Author(s), under exclusive license to Springer Nature Switzerland AG 2022
S. Avidan et al. (Eds.): ECCV 2022, LNCS 13678, pp. 376–392, 2022. https://doi.org/10.1007/978-3-031-19797-0_22

_2

Towards Interpretable Video

Super-Resolution via Alternating

Optimization

Jie Zhang¹, Jingyun Liang¹, Kai Zhang^{1(B)}, Wenguan Wang¹, Qin Wang¹,

Yulun Zhang¹, Hao Tang¹, and Luc Van Gool^{1,2}

¹Computer Vision Lab, ETH Zürich, Zürich, Switzerland

{jie.zhang, jingyun.liang, kai.zhang, wenguan.wang, qin.wang,

yulun.zhang, hao.tang, vangool}@vision.ee.ethz.ch

²KU Leuven, Leuven, Belgium

<https://github.com/caojiezhong/DAVSR>

Abstract. In this paper, we study a practical space-time video super-resolution (STVSR) problem which aims at generating a high-frame rate high-resolution sharp video from a low-frame rate low-resolution blurry video. Such problem often occurs when recording a fast dynamic event with a low-frame rate and low-resolution camera, and the captured video would suffer from three typical issues: i) motion blur occurs due to object/camera motions during exposure time; ii) motion aliasing is unavoidable when the event temporal frequency exceeds the Nyquist limit of temporal sampling; iii) high-frequency details are lost because of the low spatial sampling rate. These issues can be alleviated by a cascade of three separate sub-tasks, including video deblurring, frame interpolation, and super-resolution, which, however, would fail to capture the spatial and temporal correlations among video sequences. To address this, we propose an interpretable STVSR framework by leveraging both model-based and learning-based methods. Specifically, we formulate STVSR as a joint video deblurring, frame interpolation, and super-resolution problem, and solve it as two sub-problems in an alternate way. For the first sub-problem, we derive an interpretable analytical solution and use it as a Fourier data transform layer. Then, we propose a recurrent video enhancement layer for the second sub-problem to further recover high-frequency details. Extensive experiments demonstrate the superiority of our method in terms of quantitative metrics and visual quality.

Keywords: Video super-resolution

•Motion blur •Motion aliasing

1

Event-Based Fusion for Motion

Deblurring with Cross-modal Attention

Lei Sun^{1,2}, Christos Sakaridis², Jingyun Liang², Qijia Wang¹, Kailun Yang³,

Peng Sun¹, Yaozu Ye¹, Kaiwei Wang^{1(B)}, and Luc Van Gool^{2,4}
¹Zhejiang University, Hangzhou, China
wangkaiwei@zju.edu.cn
²ETH Zurich, Zurich, Switzerland
³KIT, Karlsruhe, Germany
⁴KU Leuven, Leuven, Belgium

Abstract. Traditional frame-based cameras inevitably suffer from motion blur due to long exposure times. As a kind of bio-inspired camera, the event camera records the intensity changes in an asynchronous way with high temporal resolution, providing valid image degradation information within the exposure time. In this paper, we rethink the event-based image deblurring problem and unfold it into an end-to-end two-stage image restoration network. To effectively fuse event and image features, we design an event-image cross-modal attention module applied at multiple levels of our network, which allows to focus on relevant features from the event branch and filter out noise. We also introduce a novel symmetric cumulative event representation specifically for image deblurring as well as an event mask gated connection between the two stages of our network which helps avoid information loss. At the dataset level, to foster event-based motion deblurring and to facilitate evaluation on challenging real-world images, we introduce the Real Event Blur (REBlur) dataset, captured with an event camera in an illumination-controlled optical laboratory. Our Event Fusion Network (EFNet) sets the new state of the art in motion deblurring, surpassing both the prior best-performing image-based method and all event-based methods with public implementations on the GoPro dataset (by up to 2.47 dB) and on our REBlur dataset, even in extreme blurry conditions. The code and our REBlur dataset are available at <https://ahupujr.github.io/EFNet/>.

1

Fast and High Quality Image Denoising
via Malleable Convolution

Yifan Jiang^{1(B)}, Bartlomiej Wronski², Ben Mildenhall², Jonathan T. Barron²,
Zhangyang Wang¹, and Tianfan Xue²

¹University of Texas at Austin, Austin, USA

yifanjiang97@utexas.edu

²Google Research, San Francisco, USA

Abstract. Most image denoising networks apply a single set of static convolutional kernels across the entire input image. This is sub-optimal for natural images, as they often consist of heterogeneous visual patterns. Dynamic convolution tries to address this issue by using per-pixel convolution kernels, but this greatly increases computational cost. In this work, we present Malleable Convolution (MalleConv), which per-

forms spatial-varying processing with minimal computational overhead. MalleConv uses a smaller set of spatially-varying convolution kernels, a compromise between static and per-pixel convolution kernels. These spatially-varying kernels are produced by an efficient predictor network running on a downsampled input, making them much more efficient to compute than per-pixel kernels produced by a full-resolution image, and also enlarging the network's receptive field compared with static kernels. These kernels are then jointly upsampled and applied to a full-resolution feature map through an efficient on-the-fly slicing operator with minimum memory overhead. To demonstrate the effectiveness of MalleConv, we use it to build an efficient denoising network we call MalleNet. MalleNet achieves high-quality results without very deep architectures, making it 8.9× faster than the best performing denoising algorithms while achieving similar visual quality. We also show that a single MalleConv layer added to a standard convolution-based backbone can significantly reduce the computational cost or boost image quality at a similar

cost. More information are on our project page: <https://yifanjiang.net/MalleConv.html> .

Keywords: Image denoising
·Dynamic kernel ·Efficiency

1

TAPE: Task-Agnostic Prior Embedding
for Image Restoration

Lin Liul, Lingxi Xie3, Xiaopeng Zhang3, Shanxin Yuan4, Xiangyu Chen5,6,
Wengang Zhou1,2, Houqiang Li1,2, and Qi Tian3(B)

1CAS Key Laboratory of Technology in GIPAS, EEIS Department,
University of Science and Technology of China, Hefei, China

2Institute of Artificial Intelligence, Hefei Comprehensive National Science Center,
Hefei, China

3Huawei Cloud BU, Shenzhen, China
tian.qil@huawei.com

4Huawei Noah's Ark Lab, London, UK
5University of Macau, Zhuhai, China

6Shenzhen Institutes of Advanced Technology, CAS, Shenzhen, China

Abstract. Learning a generalized prior for natural image restoration is an important yet challenging task. Early methods mostly involved hand-crafted priors including normalized sparsity, ℓ_0 gradients, dark channel priors, etc. Recently, deep neural networks have been used to learn various image priors but do not guarantee to generalize. In this paper, we propose a novel approach that embeds a task-agnostic prior into a transformer. Our approach, named Task-Agnostic Prior Embedding(TAPE), consists of two stages, namely, task-agnostic pre-training and task-specific fine-tuning, where the first stage embeds prior knowledge about natural images into the transformer and the second stage extracts the knowledge to assist downstream image restoration. Experiments on various types of degradation validate the effectiveness of TAPE. The image restoration performance in terms of PSNR is improved by as much as 1.45 dB and even outperforms task-specific algorithms. More importantly, TAPE shows the ability of disentangling generalized image priors from degraded images, which enjoys favorable transfer ability to unknown downstream tasks.

1

Uncertainty Inspired Underwater Image
Enhancement

Zhenqi Ful, Wuanqi, Yue Huangl, Xinghao Dingl(B), and Kai-Kuang Ma2
1Xiamen University, Fujian 361005, China

{fuzhenqi, 23320170155546}@stu.xmu.edu.cn, {yhuang2010, dxh}@xmu.edu.cn
2Nanyang Technological University, Singapore 639798, Singapore

ekkma@ntu.edu.sg

Abstract. A main challenge faced in the deep learning-based Underwater Image Enhancement (UIE) is that the ground truth high-quality image is unavailable. Most of the existing methods first generate approximate reference maps and then train an enhancement network with certainty. This kind of method fails to handle the ambiguity of the reference map. In this paper, we resolve UIE into distribution estimation and consensus process. We present a novel probabilistic network to learn the enhancement distribution of degraded underwater images. Specifically, we combine conditional variational autoencoder with adaptive instance normalization to construct the enhancement distribution. After that, we adopt a consensus process to predict a deterministic result based on a set of samples from the distribution. By learning the enhancement distribution, our method can cope with the bias introduced in the reference map labeling to some extent. Additionally, the consensus process is useful to c

capture a robust and stable result. We examined the proposed method on two widely used real-world underwater image enhancement datasets. Experimental results demonstrate that our approach enables sampling possible enhancement predictions. Meanwhile, the consensus estimate yields competitive performance compared with state-of-the-art UIE methods. Code available at <https://github.com/zhenqifu/PUIE-Net>.

Keywords: Underwater image enhancement

•Deep learning •

Probabilistic network •Adaptive instance normalization •Conditional variational autoencoder

1

Hourglass Attention Network for Image

Inpainting

Ye Deng¹, Siqi Hu¹, Rongye Meng¹, Sanping Zhou^{1,2}, and Jinjun Wang^{1(B)}

¹Xi'an Jiaotong University, Xi'an, China

{dengye, hui siqi}@stu.xjtu.edu.cn, spzhou@xjtu.edu.cn,

jinjun@mail.xjtu.edu.cn

²Shunan Academy of Artificial Intelligence, Ningbo, China

Abstract. Benefiting from the powerful ability of convolutional neural networks (CNNs) to learn semantic information and texture patterns of images, learning-based image inpainting methods have made noticeable breakthroughs over the years. However, certain inherent defects (e.g. local prior, spatially sharing parameters) of CNNs limit their performance when encountering broken images mixed with invalid information. Compared to convolution, attention has a lower inductive bias, and the output is highly correlated with the input, making it more suitable for processing images with various breakage. Inspired by this, in this paper we propose a novel attention-based network (transformer), called hourglass attention network (HAN) for image inpainting, which builds an hourglass-shaped attention structure to generate appropriate features for complemented images. In addition, we design a novel attention called Laplace attention, which introduces a Laplace distance prior for the vanilla multi-head attention, allowing the feature matching process to consider not only the similarity of features themselves, but also distance between features. With the synergy of hourglass attention structure and Laplace attention, our HAN is able to make full use of hierarchical features to mine effective information for broken images. Experiments on several benchmark datasets demonstrate superior performance by our proposed approach. The code can be found at github.com/dengye/hourglassattention.

Keywords: Image inpainting

•Attention •Transformer

1

Unfolded Deep Kernel Estimation

for Blind Image Super-Resolution

Hongyi Zheng¹, Hongwei Yong¹, and Lei Zhang^{1,2(B)}

¹The Hong Kong Polytechnic University, Hong Kong, People's Republic of China

{cshzheng, cshyong, cslzhang}@comp.polyu.edu.hk

²OPPO Research, Shenzhen, People's Republic of China

Abstract. Blind image super-resolution (BISR) aims to reconstruct a high-resolution image from its low-resolution counterpart degraded by unknown blur kernel and noise. Many deep neural network based methods have been proposed to tackle this challenging problem without considering the image degradation model. However, they largely rely on the training sets and often fail to handle images with unseen blur kernels during inference. Deep unfolding methods have also been proposed to perform BISR by utilizing the degradation model. Nonetheless, the existing deep

unfolding methods cannot explicitly solve the data term of the unfolding objective function, limiting their capability in blur kernel estimation. In this work, we propose a novel unfolded deep kernel estimation (UDKE) method, which, for the first time to our best knowledge, explicitly solves the data term with high efficiency. The UDKE based BISR method can jointly learn image and kernel priors in an end-to-end manner, and it can effectively exploit the information in both training data and image degradation model. Experiments on benchmark datasets and real-world data demonstrate that the proposed UDKE method could well predict complex unseen non-Gaussian blur kernels in inference, achieving significantly better BISR performance than state-of-the-art. The source code of UDKE is available at <https://github.com/natezhenghy/UDKE>.

Keywords: Blind image super-resolution

•Blur kernel estimation •

Unfolding method

1

Event-guided Deblurring of Unknown

Exposure Time Videos

Taewoo Kim¹, Jeongmin Lee¹, Lin Wang², and Kuk-Jin Yoon^{1(B)}

¹Korea Advanced Institute of Science and Technology, Daejeon, South Korea

{intelpro,jeanmichel,kjyoon}@kaist.ac.kr

²AI Thrust, HKUST Guangzhou and Department of CSE, HKUST,

Hong Kong, China

linwang@ust.hk

Abstract. Motion deblurring is a highly ill-posed problem due to the loss of motion information in the blur degradation process. Since event cameras can capture apparent motion with a high temporal resolution, several attempts have explored the potential of events for guiding deblurring. These methods generally assume that the exposure time is the same as the reciprocal of the video frame rate. However, this is not true in real situations, and the exposure time might be unknown and dynamically varies depending on the video shooting environment (e.g., illumination condition). In this paper, we address the event-guided motion deblurring assuming dynamically variable unknown exposure time of the frame-based camera. To this end, we first derive a new formulation for event-guided motion deblurring by considering the exposure and readout time in the video frame acquisition process. We then propose a novel end-to-end learning framework for event-guided motion deblurring. In particular, we design a novel Exposure Time-based Event Selection (ETES) module to selectively use event features by estimating the cross-modal correlation between the features from blurred frames and the events. Moreover, we propose a feature fusion module to fuse the selected features from events and blur frames effectively. We conduct extensive experiments on various datasets and demonstrate that our method achieves state-of-the-art performance. Our project code and dataset are available at: <https://intelpro.github.io/UEVD/>

1

ReCoNet: Recurrent Correction Network

for Fast and Efficient Multi-modality

Image Fusion

Zhanbo Huang¹, Jinyuan Liu², Xin Fan^{1(B)}, Risheng Liu^{1,3}, Wei Zhong¹, and Zhongxuan Luo¹

¹DUT-RU International School of Information Science and Engineering, Dalian University of Technology, Dalian, China

{xin.fan,rsliu,zhongwei,zxluo}@dlut.edu.cn

²School of Software Technology, Dalian University of Technology, Dalian, China

³Peng Cheng Laboratory, Shenzhen, China

Abstract. Recent advances in deep networks have gained great atten-

tion in infrared and visible image fusion (IVIF). Nevertheless, most existing methods are incapable of dealing with slight misalignment on source images and suffer from high computational and spatial expenses. This paper tackles these two critical issues rarely touched in the community by developing a recurrent correction network for robust and efficient fusion, namely ReCoNet. Concretely, we design a deformation module to explicitly compensate geometrical distortions and an attention mechanism to mitigate ghosting-like artifacts, respectively. Meanwhile, the network consists of a parallel dilated convolutional layer and runs in a recurrent fashion, significantly reducing both spatial and computational complexities. ReCoNet can effectively and efficiently alleviate both structural distortions and textural artifacts brought by slight misalignment. Extensive experiments on two public datasets demonstrate the superior accuracy and efficacy of our ReCoNet against the state-of-the-art IVIF methods. Consequently, we obtain a 16% relative improvement of CC on datasets with misalignment and boost the efficiency by 86%. The source code is available at <https://github.com/dlut-dimt/reconet>.

Keywords: Deep learning

· Multi-modality image fusion

1

Content Adaptive Latents and Decoder
for Neural Image Compression

Guanbo Pan¹, Guo Liu², Zhihao Hu¹, and Dong Xu^{3(B)}

¹School of Software, Beihang University, Beijing, China

²School of Computer Science and Technology, Beijing Institute of Technology, Beijing, China

³Department of Computer Science, The University of Hong Kong, Hong Kong, China
dongxu@cs.hku.hk

Abstract. In recent years, neural image compression (NIC) algorithms have shown powerful coding performance. However, most of them are not adaptive to the image content. Although several content adaptive methods have been proposed by updating the encoder-side components, the adaptability of both latents and the decoder is not well exploited. In this work, we propose a new NIC framework that improves the content adaptability on both latents and the decoder. Specifically, to remove redundancy in the latents, our content adaptive channel dropping (CACD) method automatically selects the optimal quality levels for the latents spatially and drops the redundant channels. Additionally, we propose the content adaptive feature transformation (CAFT) method to improve decoder-side content adaptability by extracting the characteristic information of the image content, which is then used to transform the features in the decoder side. Experimental results demonstrate that our proposed methods with the encoder-side updating algorithm achieve the state-of-the-art performance.

Keywords: Neural image compression

· Content adaptive coding

1

Efficient and Degradation-Adaptive
Network for Real-World Image
Super-Resolution

Jie Liang^{1,2}, Hui Zeng², and Lei Zhang^{1,2(B)}

¹The HongKong Polytechnic University, Hung Hom, Hong Kong

²OPPO Research, Shenzhen, China

csli Zhang@comp.polyu.edu.hk

Abstract. Efficient and effective real-world image super-resolution (Real-ISR) is a challenging task due to the unknown complex degradation of real-world images and the limited computation resources in practical applications. Recent research on Real-ISR has achieved significant progress

by modeling the image degradation space; however, these methods largely rely on heavy backbone networks and they are inflexible to handle images of different degradation levels. In this paper, we propose an efficient and effective degradation-adaptive super-resolution (DASR) network, whose parameters are adaptively specified by estimating the degradation of each input image. Specifically, a tiny regression network is employed to predict the degradation parameters of the input image, while several convolutional experts with the same topology are jointly optimized to specify the network parameters via a non-linear mixture of experts. The joint optimization of multiple experts and the degradation-adaptive pipeline significantly extend the model capacity to handle degradations of various levels, while the inference remains efficient since only one adaptively specified network is used for super-resolving the input image. Our extensive experiments demonstrate that DASR is not only much more effective than existing methods on handling real-world images with different degradation levels but also efficient for easy deployment. Codes, models and datasets are available at <https://github.com/csjiang/DASR>.

Keywords: Real-world image super-resolution

Degradation-adaptive Efficient super-resolution

1

Unidirectional Video Denoising by
Mimicking Backward Recurrent Modules
with Look-Ahead Forward Ones

Junyi Lil, Xiaohe Wul(B), Zhenxing Niu2, and Wangmeng Zuol

1Harbin Institute of Technology, Harbin, China

csxhwu@gmail.com ,wmzuo@hit.edu.cn

2Xidian University, Xi'an, China

Abstract. While significant progress has been made in deep video denoising, it remains very challenging for exploiting historical and future frames. Bidirectional recurrent networks (BiRNN) have exhibited appealing performance in several video restoration tasks. However, BiRNN is intrinsically offline because it uses backward recurrent modules to propagate from the last to current frames, which causes high latency and large memory consumption. To address the offline issue of BiRNN, we present a novel recurrent network consisting of forward and look-ahead recurrent modules for unidirectional video denoising. Particularly, look-ahead module is an elaborate forward module for leveraging information from near-future frames. When denoising the current frame, the hidden features by forward and look-ahead recurrent modules are combined, thereby making it feasible to exploit both historical and near-future frames. Due to the scene motion between non-neighboring frames, border pixels missing may occur when warping look-ahead feature from near-future frame to current frame, which can be largely alleviated by incorporating forward warping and proposed border enlargement. Experiments show that our method achieves state-of-the-art performance with constant latency and memory consumption. Code is available at <https://github.com/nagejacob/FloRNN>.

Keywords: Video denoising

Recurrent neural networks Temporal alignment

1

Self-supervised Learning for Real-World
Super-Resolution from Dual Zoomed
Observations

Zhilu Zhangl, Ruohao Wangl, Hongzhi Zhangl(B), Yunjin Chen1,2,
and Wangmeng Zuol,2

1Harbin Institute of Technology, Harbin, China

zhanghz0451@gmail.com, wmzuo@hit.edu.cn

2Peng Cheng Laboratory, Shenzhen, China

Abstract. In this paper, we consider two challenging issues in reference-based super-resolution (RefSR), (i) how to choose a proper reference image, and (ii) how to learn real-world RefSR in a self-supervised manner. Particularly, we present a novel self-supervised learning approach for real-world image SR from observations at dual camera zooms (SelfDZSR). Considering the popularity of multiple cameras in modern smartphones, the more zoomed (telephoto) image can be naturally leveraged as the reference to guide the SR of the lesser zoomed (short-focus) image. Furthermore, SelfDZSR learns a deep network to obtain the SR result of short-focus image to have the same resolution as the telephoto image. For this purpose, we take the telephoto image instead of an additional high-resolution image as the supervision information and select a center patch from it as the reference to super-resolve the corresponding short-focus image patch. To mitigate the effect of the misalignment between short-focus low-resolution (LR) image and telephoto ground-truth (GT) image, we design an auxiliary-LR generator and map the GT to an auxiliary-LR while keeping the spatial position unchanged. Then the auxiliary-LR can be utilized to deform the LR features by the proposed adaptive spatial transformer networks (AdaSTN), and match the Ref features to GT. During testing, SelfDZSR can be directly deployed to super-solve the whole short-focus image with the reference of telephoto image. Experiments show that our method achieves better quantitative and qualitative performance against state-of-the-arts. Codes are available at <https://github.com/cszhilul998/SelfDZSR>.

Keywords: Reference-based super-resolution

· Self-supervised

learning · Real world

1

Secrets of Event-Based Optical Flow

Shintaro Shibata^{1,2(B)}, Yoshimitsu Aoki¹, and Guillermo Gallego^{2,3}

¹Department of Electronics and Electrical Engineering, Faculty of Science and Technology, Keio University, Kanagawa, Japan

sshiba@keio.jp

²Department of EECS, Technische Universität Berlin, Berlin, Germany

³Einstein Center Digital Future and SCIOI Excellence Cluster, Berlin, Germany

Abstract. Event cameras respond to scene dynamics and offer advantages to estimate motion. Following recent image-based deep-learning achievements, optical flow estimation methods for event cameras have rushed to combine those image-based methods with event data. However, it requires several adaptations (data conversion, loss function, etc.) as they have very different properties. We develop a principled method to extend the Contrast Maximization framework to estimate optical flow from events alone. We investigate key elements: how to design the objective function to prevent overfitting, how to warp events to deal better with occlusions, and how to improve convergence with multi-scale raw events. With these key elements, our method ranks first among unsupervised methods on the MVSEC benchmark, and is competitive on the DSEC benchmark. Moreover, our method allows us to expose the issues of the ground truth flow in those benchmarks, and produces remarkable results when it is transferred to unsupervised learning settings. Our code is available at <https://github.com/tub-rip/event-based-optical-flow>.

1

Towards Efficient and Scale-Robust

Ultra-High-Definition Image Demoiréing

Xin Yu¹, Peng Dai¹, Wenbo Li², Lan Ma³, Jiajun Shen³, Jialin Li⁴, and Xiaojuan Qil^(B)

¹The University of Hong Kong, Hong Kong, China

xjq@eee.hku.hk

²The Chinese University of Hong Kong, Hong Kong, China

³TCL AI Lab, Hong Kong, China

⁴Sun Yat-sen University, Guangzhou, China

Abstract. With the rapid development of mobile devices, modern widely-used mobile phones typically allow users to capture 4K resolution (i.e., ultra-high-definition) images. However, for image demoiréing, a challenging task in low-level vision, existing works are generally carried out on low-resolution or synthetic images. Hence, the effectiveness of these methods on 4K real-world images is still unknown. In this paper, we explore

moiré pattern removal for ultra-high-definition images. To this end, we propose the first ultra-high-definition demoiréing dataset (UHDM), which contains 5,000 real-world 4K resolution image pairs, and conduct a benchmark study on current state-of-the-art methods. Further, we present an efficient baseline model ESDNet for tackling 4K moiré images, wherein we build a semantic-aligned scale-aware module to address the scale variation of moiré patterns. Extensive experiments manifest the effectiveness of our approach, which outperforms state-of-the-art methods by a large margin while being much more lightweight. Code and dataset are available at <https://xinyu-andy.github.io/uhdm-page>.

Keywords: Image demoiréing

·Image restoration ·

Ultra-high-definition

1

ERDN: Equivalent Receptive Field

Deformable Network for Video

Deblurring

Bangrui Jiang^{1,2}, Zhihuai Xie^{2(B)}, Zhen Xia², Songnan Li², and Shan Liu²

¹Tsinghua Shenzhen International Graduate School,

Tsinghua University, Beijing, China

²Tencent Media Lab, Shenzhen, China

{zhihuaixie,zhenxia,sunnysnli,shanli}@tencent.com

Abstract. Video deblurring aims to restore sharp frames from blurry video sequences. Existing methods usually adopt optical flow to compensate misalignment between reference frame and each neighboring frame. However, in accurate flow estimation caused by large displacements will lead to artifacts in the warped frames. In this work, we propose an equivalent receptive field deformable network (ERDN) to perform alignment at the feature level without estimating optical flow. The ERDN introduces a dual pyramid alignment module, in which a feature pyramid is constructed to align frames using deformable convolution in a cascaded manner. Specifically, we adopt dilated spatial pyramid blocks to predict offsets for deformable convolutions, so that the theoretical receptive field is equivalent for each feature pyramid layer. To restore the sharp frame, we propose a gradient guided fusion module, which incorporates structure priors into the restoration process. Experimental results demonstrate that the proposed method outperforms previous state-of-the-art methods on multiple benchmark datasets. The code is made available at:

<https://github.com/TencentCloud/ERDN>.

Keywords: Video deblurring

·Deformable convolution ·Receptive

field

1

Rethinking Generic Camera Models

for Deep Single Image Camera Calibration

to Recover Rotation and Fisheye

Distortion

Nobuhiko Wakai¹(B), Satoshi Sato¹, Yasunori Ishii¹,
and Takayoshi Yamashita²

¹Panasonic Holdings, Osaka, Japan

{wakai.nobuhiko,sato.satoshi,ishii.yasunori}@jp.panasonic.com

²Chubu University, Aichi, Japan

{takayoshi@isc.chubu.ac.jp}

Abstract. Although recent learning-based calibration methods can predict extrinsic and intrinsic camera parameters from a single image, the accuracy of these methods is degraded in fisheye images. This degradation is caused by mismatching between the actual projection and expected projection. To address this problem, we propose a generic camera model that has the potential to address various types of distortion. Our generic camera model is utilized for learning-based methods through a closed-form numerical calculation of the camera projection.

Simultaneously to recover rotation and fisheye distortion, we propose a learning-based calibration method that uses the camera model. Furthermore, we propose a loss function that alleviates the bias of the magnitude of errors for four extrinsic and intrinsic camera parameters. Extensive experiments demonstrated that our proposed method outperformed conventional methods on two large-scale datasets and images captured by off-the-shelf fisheye cameras. Moreover, we are the first researchers to analyze the performance of learning-based methods using various types of projection for off-the-shelf cameras.

Keywords: Camera calibration

·Fisheye camera ·Rectification

1

ART-SS: An Adaptive Rejection

Technique for Semi-supervised

Restoration for Adverse

Weather-Affected Images

Rajeev Yasarla¹(B), Carey E. Priebe, and Vishal M. Patel

Johns Hopkins University, Baltimore, MD 21218, USA

{ryasar11,cep,vpatel36}@jhu.edu

Abstract. In recent years, convolutional neural network-based single image adverse weather removal methods have achieved significant performance improvements on many benchmark datasets. However, these methods require large amounts of clean-weather degraded image pairs for training, which is often difficult to obtain in practice. Although various weather degradation synthesis methods exist in the literature, the use of synthetically generated weather degraded images often results in sub-optimal performance on the real weather degraded images due to the domain gap

between synthetic and real world images. To deal with this problem, various semi-supervised restoration (SSR) methods have been proposed for deraining or dehazing which learn to restore clean image using synthetically generated datasets while generalizing better using unlabeled real-world images. The performance of a semi-supervised method is essentially based on the quality of the unlabeled data. In particular, if the unlabeled data characteristics are very different from that of the labeled data, then the performance of a semi-supervised method degrades significantly. We theoretically study the effect of unlabeled data on the performance of an SSR method and develop a technique that rejects the unlabeled images that degrade the performance. Extensive experiments and ablation study show that the proposed sample rejection method increases the performance of existing SSR deraining and dehazing methods significantly. Code is available at: <https://github.com/rajeevyasarla/ART-SS>.

Keywords: Semi-supervision

·Deraining ·Dehazing ·Rejection

technique

Fusion from Decomposition:

A Self-Supervised Decomposition

Approach for Image Fusion

Pengwei Liang¹, Junjun Jiang^{1(B)}, Xianming Liu¹, and Jiayi Ma²

¹Harbin Institute of Technology, Harbin 150001, China

{jiangjunjun, csxm}@hit.edu.cn

²Wuhan University, Wuhan 430072, China

Abstract. Image fusion is famous as an alternative solution to generate one high-quality image from multiple images in addition to image restoration from a single degraded image. The essence of image fusion is to integrate complementary information or best parts from source images.

The current fusion methods usually need a large number of paired samples or sophisticated loss functions and fusion rules to train the supervised or unsupervised model. In this paper, we propose a powerful image decomposition model for fusion task via the self-supervised representation learning, dubbed Decomposition for Fusion (DeFusion). Without any paired data or sophisticated loss, DeFusion can decompose the source images into a feature embedding space, where the common and unique features can be separated. Therefore, the image fusion can be achieved within the embedding space through the jointly trained reconstruction (projection) head in the decomposition stage even without any

fine-tuning. Thanks to the development of self-supervised learning, we can train the model to learn image decomposition ability with a brute

but simple pretext task. The pretrained model allows for learning very effective features that generalize well: the DeFusion is a unified versatile

framework that is trained with an image fusion irrelevant dataset and can be directly applied to various image fusion tasks. Extensive experiments demonstrate that the proposed DeFusion can achieve comparable or even better performance compared to state-of-the-art methods (whether supervised or unsupervised) for different image fusion tasks.

Keywords: Image fusion

• Self-supervised learning • Image decomposition

1

Learning Degradation Representations

for Image Deblurring

Dasong Li¹, Yi Zhang¹, Ka Chun Cheung², Xiaogang Wang^{1,4},

Hongwei Qin^{3(B)}, and Hongsheng Li^{1,4,5(B)}

¹MMLab, CUHK, Hong Kong, China

dasongli@link.cuhk.edu.hk, hsli@ee.cuhk.edu.hk

²NVIDIA AI Technology Center, Santa Clara, USA

³SenseTime Research, Hong Kong, China

qinhongwei@sensetime.com

⁴Centre for Perceptual and Interactive Intelligence Limited, Hong Kong, China

⁵Xidian University, Xi'an, China

Abstract. In various learning-based image restoration tasks, such as image denoising and image super-resolution, the degradation representations were widely used to model the degradation process and handle complicated degradation patterns. However, they are less explored in learning-based image deblurring as blur kernel estimation cannot perform well in real-world challenging cases. We argue that it is particularly necessary for image deblurring to model degradation representations since blurry patterns typically show much larger variations than noisy patterns or high-frequency textures. In this paper, we propose a framework to learn spatially adaptive degradation representations of blurry images. A novel joint image reblurring and deblurring learning process is presented to improve the expressiveness of degradation rep-

representations. To make learned degradation representations effective in reblurring and deblurring, we propose a Multi-Scale Degradation Injection Network (MSDI-Net) to integrate them into the neural networks. With the integration, MSDI-Net can handle various and complicated blurry patterns adaptively. Experiments on the GoPro and RealBlur datasets demonstrate that our proposed deblurring framework with the learned degradation representations outperforms state-of-the-art methods with appealing improvements. The code is released at <https://github.com/dasonglil/Learning-degradation>.

Keywords: Image deblurring · Degradation representations

Supplementary Information The online version contains supplementary material available at https://doi.org/10.1007/978-3-031-19797-0_42.

© The Author(s), under exclusive license to Springer Nature Switzerland AG 2022
S. Avidan et al. (Eds.): ECCV 2022, LNCS 13678, pp. 736–753, 2022. <https://doi.org/10.1007/978-3-031-19797-0>

_4
