

# Deep Burst Super-Resolution

Goutam Bhat      Martin Danelljan      Luc Van Gool      Radu Timofte  
Computer Vision Lab, ETH Zurich, Switzerland

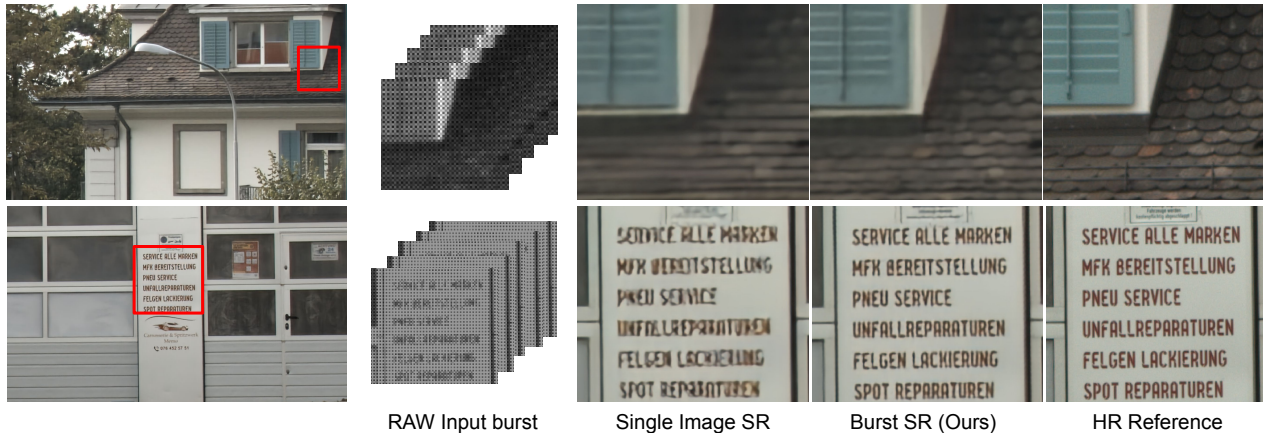


Figure 1. Our network generates a super-resolved RGB image from an input burst consisting of multiple noisy RAW frames. In contrast to the single image baseline, our approach combines information from multiple frames to obtain a more detailed reconstruction of the scene. The results shown are for super-resolution by a factor of 4.

## Abstract

While single-image super-resolution (SISR) has attracted substantial interest in recent years, the proposed approaches are limited to learning image priors in order to add high frequency details. In contrast, multi-frame super-resolution (MFSR) offers the possibility of reconstructing rich details by combining signal information from multiple shifted images. This key advantage, along with the increasing popularity of burst photography, have made MFSR an important problem for real-world applications.

We propose a novel architecture for the burst super-resolution task. Our network takes multiple noisy RAW images as input, and generates a denoised, super-resolved RGB image as output. This is achieved by explicitly aligning deep embeddings of the input frames using pixel-wise optical flow. The information from all frames are then adaptively merged using an attention-based fusion module. In order to enable training and evaluation on real-world data, we additionally introduce the BurstSR dataset, consisting of smartphone bursts and high-resolution DSLR ground-truth. We perform comprehensive experimental analysis, demonstrating the effectiveness of the proposed architecture.

## 1. Introduction

Super-resolution (SR) is the task of generating a high-resolution (HR) image, given one or several low-resolution (LR) observations. It is a widely studied problem [6, 8, 20, 21, 23, 24, 26, 29, 39, 42, 45, 48, 51] with numerous practical applications. In recent years, the SR community has mainly focused on the single image super-resolution (SISR) task, where an HR image is estimated from a single LR input. Due to the ill-posed nature of the SISR problem, these methods are limited to adding high frequency details through learned image priors.

The multi-frame super-resolution (MFSR), on the other hand, aims to reconstruct the original HR image using multiple LR images. If the input images have sub-pixel shifts with respect to each other, due to *e.g.* camera motion, they provide different LR samplings of the underlying scene. MFSR approaches can thus exploit this additional signal information to generate a higher quality image, compared to the SISR approaches (see Fig. 1). The MFSR problem naturally arises in the increasingly popular mobile burst photography, where the images have different sub-pixel shifts due to natural hand tremors [45]. This opens up the possibility of using MFSR to overcome the resolution constraint

in mobile cameras imposed by the cost and size restrictions.

Despite the aforementioned advantages, MFSR has received little attention in recent years. This is in stark contrast to SISR, where deep learning has led to significant advancements in SR performance. Compared to the SISR case, the MFSR problem imposes significant challenges when developing deep learning based solutions. Firstly, a MFSR architecture must be able to align the noisy input frames with sub-pixel accuracy in order to enable fusion. Secondly, it should be able to effectively fuse the information from the aligned frames, while being robust to alignment errors. Furthermore, the lack of benchmark datasets for the general MFSR task has led to a limited interest in the MFSR problem. We address these issues by proposing a novel deep learning based approach for the MFSR problem, along with a real-world dataset.

Our network directly operates on noisy RAW bursts captured from a hand-held camera and generates a denoised, demosaicked, and super-resolved image as output. This is achieved by developing a novel attention-based fusion module which can adaptively merge an arbitrary number of input frames in order to produce a high quality output. Our approach is not limited to simple motions between the images, such as translation or homography. Instead, we estimate dense pixel-wise optical flow to align the deep feature encoding of each input frame. The aligned representations of each frame are then merged by computing element-wise fusion weights. This allows the network to adaptively select the reliable and informative content from each image, while discarding, *e.g.*, misaligned regions.

The conventional approach in SISR is to train and evaluate models on synthetically generated data. However, this has been shown to not generalize to real-world images due to inaccuracies in data generation model [3, 30, 31, 32]. Accurately modelling the image formation process for MFSR is further challenging due to the additional complexity introduced by camera motion. We therefore introduce the BurstSR dataset: the first real-world burst super-resolution dataset. Our dataset consists of 200 RAW bursts captured using a hand held mobile camera. Furthermore, we provide a high quality HR ground truth for each burst using a DSLR with zoom lens. We believe that our BurstSR dataset can serve as a valuable benchmark and source of training data to stimulate future research in MFSR.

**Contributions:** Our main contributions are summarized as follows. **(i)** We introduce the first real world burst super-resolution dataset consisting of RAW bursts and corresponding HR ground truths. **(ii)** We propose a novel MFSR architecture which can perform joint denoising, demosaicking, and SR using bursts captured from a handheld camera. **(iii)** Our architecture employs an attention-based fusion method to adaptively merge the input images to generate high quality HR output **(iv)** We further address mis-

alignment issues encountered when training on real world data by introducing a loss function which can internally correct these mis-alignments.

We perform comprehensive experiments on a synthetic dataset, as well as the BurstSR test set, in order to validate our contributions. Our approach demonstrates promising SR performance on real world bursts, significantly outperforming alternative methods in a user study. We also provide a detailed ablative study, analysing the impact of key components in the proposed MFSR architecture.

## 2. Related Work

**Single Image Super-Resolution:** SISR is a widely studied task with a variety of proposed methods, for example based on the frequency domain [19, 34, 36], interpolation techniques [5, 16, 25], sparse representations [28, 46, 47] or patch and examples [4, 11, 13]. Dong *et al.* [7] were the first to train a deep CNN to directly map the input LR image to the HR output. A number of approaches have subsequently improved upon this work using more effective network architectures [8, 21, 23, 26, 39, 51] and loss functions [20, 24, 29, 42, 48].

**Multi-Frame Super-Resolution:** Compared to SISR approaches which solely rely on image priors to perform super-resolution, MFSR methods aim to merge multiple aliased images of the same scene to reconstruct a higher resolution output. The MFSR problem was first addressed by Tsai and Huang [40], who proposed a frequency domain based method that assumes known translations between input images. Later, Peleg *et al.* [35] and Irani and Peleg [18] introduced the iterative back-projection approach. They estimate an initial HR image and simulate the imaging process to generate the LR images. The reconstruction error between the generated and input LR images is then minimized iteratively to refine the HR image. Hardie *et al.* [14] extended this approach with an improved observation model and a regularization term. Farsui *et al.* [9] proposed a joint multi-frame demosaicking and super-resolution approach using a maximum a posteriori estimation framework. Zomet *et al.* [52] use information from multiple sensors to perform super-resolution. Recently, Wronski *et al.* [45] proposed a MFSR method for hand-held cameras, where a kernel regression technique is employed to merge aligned input frames robustly. Unlike in SISR, only a few deep learning based approaches have been proposed for MFSR. Ustinova and Lempitsky [41] proposed a multi-frame network for face super-resolution. Deudon *et al.* [6] developed HighRes-net, a MFSR network for satellite imagery. HighRes-net aligns each input frame to a reference frame implicitly, and merges them using a recursive fusion method. Another approach for satellite imagery, namely DeepSUM [33], assumes only translation motion between

frames and utilizes 3D convolution for fusion. In contrast to these previous approaches which are focused on remote sensing, we tackle the general problem of burst SR from any handheld camera.

**Learning real world super-resolution:** SR approaches are commonly trained using synthetically generated LR images. However, such a training strategy has been shown to not generalize well to real-world images [3, 30, 31, 32]. A few recent works have tried to address this issue by learning real world degradation models [3, 30, 44]. Another approach is to learn camera specific SR models directly using real world data. Such a strategy allows the network to learn the characteristics of the particular sensor, leading to improved performance [50]. This is however challenging due to difficulties in collecting paired training data for SR. Zhang *et al.* [50] address this by using LR-HR pairs captured using a zoom lens for training. In order to handle the spatial and color mis-alignments between LR-HR pairs, a novel contextual bilateral loss is employed for training.

### 3. Burst Super-Resolution Network

In this section, we describe our burst super-resolution network. Our network inputs multiple noisy, RAW, low-resolution (LR) images captured in a single burst. The architecture processes and combines the information in individual images to generate a high-resolution (HR) RGB image as output. Thus, our network performs joint denoising, demosaicking, and SR. Since the images in a burst are captured in a rapid sequence from a hand-held device, they include small inter-frame offsets. This ensures multiple aliased versions of the same scene, providing additional signal information for SR. Consequently, by effectively merging the information from the whole burst, our network can better reconstruct the underlying scene to generate a higher quality output, compared to single frame approaches.

An overview of our architecture is shown in Figure 2. Our network takes a RAW burst sequence  $\{b_i\}_{i=1}^N$  of any arbitrary size  $N$  as input. Here, each image  $b_i \in \mathbb{R}^{W \times H}$  is the RAW sensor data obtained from the camera. The images in the burst are first encoded independently in order to obtain deep feature representations  $\{e_i\}_{i=1}^N$ . Next, we align and warp each of the feature maps to a common reference frame  $b_1$  using the offsets estimated by an alignment network. The aligned feature maps are then combined by our fusion module to obtain a merged feature map  $\hat{e}$ . We propose an attention-based fusion approach that predicts element-wise fusion weights. This allows the network to adaptively select the most useful information from each image in the burst. The merged feature map is then passed to the decoder module which outputs the final RGB image  $y \in \mathbb{R}^{sW \times sH \times 3}$ , where  $s$  is the super-resolution factor. We detail each network module of our architecture in the subsequent sections.

### 3.1. Encoder

The encoder module  $E$  independently maps each input burst image  $b_i$  to a deep feature representation  $e_i$ . To ensure translational invariance, we first pack each  $2 \times 2$  block in the raw Bayer pattern along the channel dimension, obtaining a 4 channel image  $\tilde{b}_i \in \mathbb{R}^{\frac{W}{2} \times \frac{H}{2} \times 4}$  at half the initial resolution. This LR image is passed through the encoder, consisting of an initial convolutional layer followed by a series of residual blocks. In order to achieve a high-dimensional encoding that allows more effective fusion of several frames, we radically expand the feature dimensionality with a final convolutional layer. The resulting  $D$ -dimensional encoding  $E(\tilde{b}_i) = e_i \in \mathbb{R}^{\frac{W}{2} \times \frac{H}{2} \times D}$  thus achieves a rich embedding of the input image. We use  $D = 512$  in our experiments.

### 3.2. Alignment Module

One of the important challenges in burst SR is that the pixel-wise displacement between the images is unknown. The displacements stem from both global camera motion and scene variations. In order to achieve an effective fusion of multiple frames, the information first needs to be aligned. We address this problem by explicitly aligning the individual image embeddings  $e_i$  to a common reference LR image, called the *base frame*. For convenience, we let the first image  $\tilde{b}_1$  denote the base frame. Camera motion is often modelled using a homography when imaging static and distant scenes. However, we found these assumptions to seldom hold in the real-world scenario. Thus, we allow greater flexibility in our alignment module by computing dense pixel-wise optical flow  $f_i \in \mathbb{R}^{\frac{W}{2} \times \frac{H}{2} \times 2}$  between every burst image  $\tilde{b}_i$  and the reference image  $\tilde{b}_1$ . Pixel-wise flow can capture global camera motion while also accounting for any object motion in the scene. The estimated flow vectors  $f_i$  are then used to warp the feature maps  $e_i$  to the base frame using a bilinear kernel

$$\tilde{e}_i = \phi(e_i, f_i), \quad f_i = F(\tilde{b}_i, \tilde{b}_1) \quad (1)$$

Here,  $\phi$  denotes the warping operation,  $F$  is the flow estimator, while  $\tilde{e}_i$  is the warped feature map. The warped feature maps  $\{\tilde{e}_i\}_{i=1}^N$ , as well as the computed flow vectors  $\{f_i\}_{i=1}^N$  are then passed to the fusion module. Here, the flow vectors  $f_1$  for the base frame is set to 0. We use a state-of-the-art optical flow network PWC-Net [38] as our flow estimator  $F$  due to its high accuracy and speed. Since PWC-Net is trained to operate on RGB images, we discard one of the two green channels in  $\tilde{b}_i$  to generate input RGB images.

### 3.3. Fusion Module

The fusion module combines information across the individual burst images to generate a merged feature embedding  $\hat{e}$ . In order to be able to operate on bursts of arbitrary sizes, the fusion module must be able to merge any number of input frames. Consequently, it is infeasible to *e.g.*

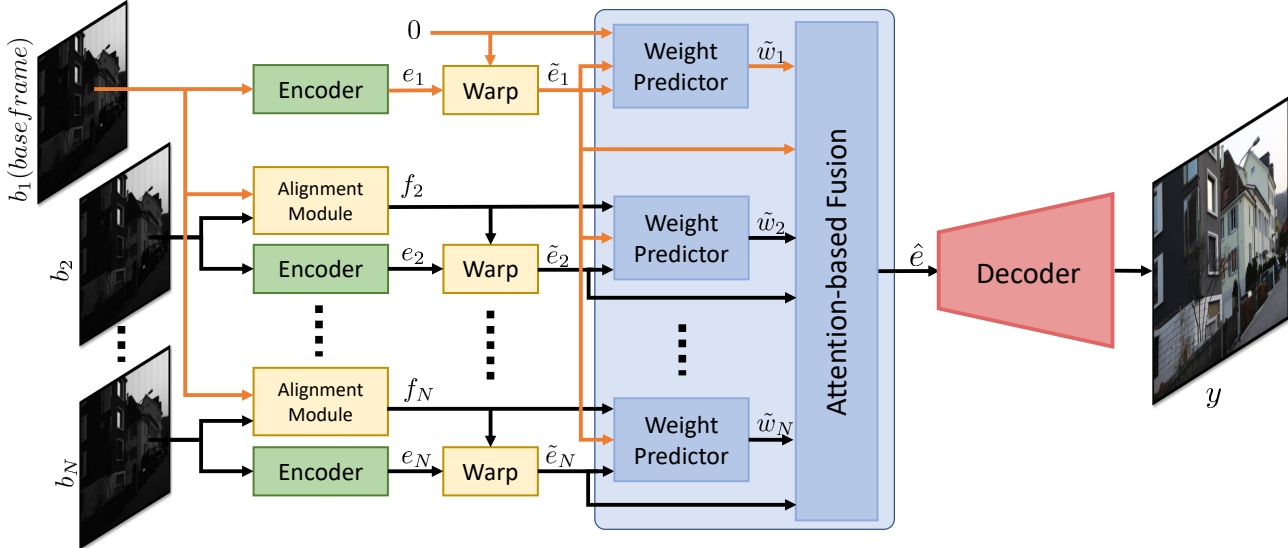


Figure 2. An overview of our burst super-resolution architecture. Each image  $b_i$  in the input burst is first passed independently through the encoder. The resulting feature maps are then warped to the base frame ( $b_1$ ) coordinates using the flow vectors  $f_i$  predicted by the alignment predictor. The aligned feature maps are then merged using an attention-based fusion module, using fusion weights computed by the weight predictor. The merged feature map  $\hat{e}$  is then passed through to decoder module to obtain the super-resolved RGB image as output.

directly concatenate the input feature maps along the channel dimension. We further found simple pooling operations such as element-wise max or average pool across the burst to provide unsatisfactory results. This is because the fusion module needs to be able to merge adaptively based on *e.g.* image content, noise levels, etc. For instance, it can be beneficial to have uniform fusion weights for textureless regions in order to perform denoising. On the other hand, it is preferable to have low fusion weights for any mis-aligned frame in order to avoid ghosting artifacts. We therefore propose an attention-based fusion approach, where element-wise fusion weights are predicted by a weight predictor network  $W$ . This provides flexibility to the network to effectively extract the useful information from each image, while also being able to process arbitrary number of input images.

The weight predictor network  $W$  utilizes both the aligned feature maps  $\tilde{e}_i$  and the flow vectors  $f_i$  to estimate the unnormalized attention weights  $\tilde{w}_i \in \mathbb{R}^{\frac{W}{2} \times \frac{H}{2} \times D}$  for each embedding  $\tilde{e}_i$ . We first project  $\tilde{e}_i$  to a lower dimension feature map  $\tilde{e}_i^p$  for computational efficiency. To compute the attention weights for  $\tilde{e}_i$ , we use the projected base frame feature map  $\tilde{e}_1^p$ , as well as the residual  $r_i = \tilde{e}_i^p - \tilde{e}_1^p$  between  $\tilde{e}_i^p$  and  $\tilde{e}_1^p$ . The base frame map  $\tilde{e}_1^p$  contains information about the local image content. This is informative to determine *e.g.* whether to use uniform fusion weights to achieve denoising, or perform edge-aware fusion in order to avoid over smoothing edges. On the other hand, the residual  $r_i$  can provide an estimate of alignment errors and thus help assign low fusion weights to misaligned regions. Additionally, we use the flow vectors  $f_i$  for weight estimation as they provide the sub-pixel sampling location of the image

data. We obtain the sub-pixel offset by computing modulo 1 of the flow vectors  $f_i$  and pass it through a small CNN to obtain the flow features  $\hat{f}_i$ . The reference frame features  $\tilde{e}_1^p$ , the feature residual  $r_i$ , and the flow features  $\hat{f}_i$  are concatenated along the channel dimension and passed through a residual network to obtain the raw fusion weights  $\tilde{w}_i$ . The raw fusion weights are then normalized across the burst using a softmax function to obtain the final attention weights  $w_i$ . The merged feature map  $\hat{e}$  is then obtained as the following weighted sum,

$$\hat{e} = \sum_{i=1}^N w_i \cdot e_i, \quad w_i = \frac{e^{\tilde{w}_i}}{\sum_j e^{\tilde{w}_j}}, \quad \tilde{w}_i = W(\tilde{e}_1, r_i, \hat{f}_i). \quad (2)$$

Here,  $\cdot$  denotes element-wise multiplication. The merged feature map  $\hat{e}$  is then passed to the decoder module to generate the final output.

### 3.4. Decoder

The decoder module generates the output high-resolution RGB image from the fused feature map  $\hat{e}$ . We first project the input feature map to 128 channels and pass it through a residual network. Next, we upsample this to the desired resolution  $sH \times sW$  using sub-pixel convolution [37]. We use a convolution layer to increase the feature dimension to  $2^2 s^2 D'$ , obtaining a tensor of shape  $\frac{H}{2} \times \frac{W}{2} \times 2^2 s^2 D'$ . The feature vectors at each spatial location are then re-arranged into a  $2s \times 2s \times D'$  map to obtain a higher resolution feature map of shape  $H \times W \times D'$ . Here,  $D'$  is the output feature dimension of the sub-pixel convolution layer. Compared to performing naïve upsampling using

*e.g.* bilinear interpolation, sub-pixel convolution allows us to effectively decode the sub-pixel information encoded in the different feature channels. In order to avoid checkerboard artifacts, we use the ICNR initialization [1] for the sub-pixel convolution layer and additionally apply Gaussian smoothing to its output. The upsampled feature map is then passed through another set of residual blocks, followed by a conv layer to obtain the high resolution RGB image  $y$ .

## 4. BurstSR Dataset

The aim of this work is to propose a burst SR method for real-world photography applications. In order to validate the performance of our approach, it is essential to train and evaluate our models on real data. Hence, we collect a new dataset, called BurstSR. To the best of our knowledge, it is the first real world burst super-resolution dataset. The BurstSR dataset consists of 200 RAW burst sequences, and corresponding high-resolution ground truths. Each burst sequence contains 14 RAW images captured using identical camera settings (*e.g.* exposure, ISO). All bursts are captured using a handheld smartphone camera. Our dataset therefore contains natural hand tremors, resulting in small random offsets between the images within a burst that are essential for MFSR [40]. For each burst sequence, we also capture a high-resolution image using a DSLR camera mounted on a tripod to serve as ground truth. Our BurstSR dataset will be released upon publication. We believe that it can serve as an important training set and benchmark for the community, in order to raise the interest in the important MFSR problem.

We capture the burst images in our dataset using a handheld Samsung Galaxy S8 smartphone camera. In order to capture and store RAW bursts, we developed a custom app using Camera2 API. On pressing the shutter, the app runs the camera’s auto-focus, auto-exposure, and auto-white-balance algorithms to determine the camera settings. These settings are then used to capture a fixed number of RAW images. The corresponding ground truth images for each burst are collected using a Canon 5D Mark IV DSLR camera mounted on a tripod. We use a zoom lens with a focal length of 70mm to obtain images with  $\approx 4$  times higher spatial resolution compared to burst images captured from the phone camera. The images are taken using a smaller aperture size (F18) to have a wider depth of field. Other capture settings are automatically determined by the camera. We hold the phone camera just above the DSLR when taking bursts in order to minimize misalignments between the two images. Additionally, we use a timer on the DSLR to synchronize the capture time between the two cameras. In order to minimize the effect of any error in temporal synchronization, we try to capture static scenes with little (*e.g.* leaves moving due to wind) or no motion. We collect 200 bursts in total, which are split into train, validation, and test sets consisting of 160, 20, and 20 sequences, respectively.

## 5. Training

In this section, we describe our training pipeline in detail. Due to the high cost and effort associated with collecting real-world paired data for MFSR, it is impractical to obtain large scale real world datasets for training our model from scratch. We therefore exploit methods for synthetic data generation to first pre-train our networks. The resulting model serves as a strong initialization, which is then finetuned on our BurstSR dataset to perform real-world SR.

### 5.1. Synthetic data training

We generate synthetic RAW bursts for pre-training our model using the sRGB images from the training split of Zurich RAW to RGB dataset [17]. Given a sRGB image, we apply the inverse camera pipeline described in [2] to obtain raw sensor values. Next, we generate a synthetic burst of size  $N$  by applying random translations and rotations to the converted RGB image. The translation and rotation values are sampled independently from the range  $[-24, 24]$  pixels and  $[-1, 1]$  degrees, respectively. The transformed images are then downsampled by the desired super-resolution factor  $s$  to obtain the low resolution RGB burst. We use bilinear kernel for both image translation/rotation and downsampling. Next, we add shot and read noise to the burst images, as described in [2]. We then discard two color channels per pixel according to the Bayer CFA to obtain the mosaicked RAW burst. We extract  $96 \times 96$  crops from the resulting RAW burst for our training. Our network is trained in a fully supervised manner by minimizing the  $L_1$  loss between the network prediction and the ground truth image. The loss is computed in the linear sensor space, before any post processing *e.g.* gamma compression or tone-mapping.

### 5.2. Real data training

In order to reconstruct the HR image using multiple aliased LR observations, a MFSR model needs to learn the image formation process in a camera. However, due to differences in the image formation process in a real camera and the one modelled by our synthetic pipeline, a network trained on only synthetic data is thus expected to have sub-optimal performance when applied to real data. Hence, we fine-tune the pre-trained synthetic data model on our real world BurstSR dataset in order to adapt the model to the particular camera sensor.

**Data Processing:** Here, we describe the pipeline used to pre-process the collected BurstSR data for training. Since the images captured using phone and DSLR cameras have different field of views (FOV), we first crop out matching field of view from each image in the burst. This is done by estimating a homography between the first image in the burst and DSLR image using SIFT [27] and RANSAC [10]. Next, we extract  $160 \times 160$  crops from the burst images in a

sliding window manner, with a stride of 80 pixels. For each crop, we again estimate homography between the crop and the corresponding region in DSLR image to perform local alignment. The aligned DSLR image region is then down-sampled to  $160s \times 160s$  to obtain the ground truth crop. In order to filter out crops with incorrect alignment, we discard phone-DSLR pairs which have a normalized cross correlation of less than 0.9 between them.

**Training loss:** There are several challenges when training our model on real bursts due to the unavoidable mis-alignments between the input burst and the ground truth. Firstly, even though we align the burst images to DSLR using homography, there can still be misalignments between the pair due to perspective shifts, error in homography estimation, etc. Secondly, since the burst and ground truth images are captured using two different sensors, there is a color mis-match between the two. Thus, it is not feasible to train the model by directly computing a pixel-wise error between the network prediction  $y$  and the ground truth  $y_{GT}$ .

In order to handle the spatial mis-alignment issue, we first estimate the optical flow  $f_{Pred,GT}$  between the prediction and ground truth using PWC-Net. The estimated flow is then used to warp the network prediction to the ground truth co-ordinates. Next, we estimate a global color mapping between the burst and the ground truth in order to handle the color mis-match. We first downsample the ground truth image to the same resolution as the input burst images. The estimated flow  $f_{Pred,GT}$  is then used to align the first image in the burst to the downsampled ground truth. In order to minimize the effect of small mis-alignments, we apply Gaussian smoothing on both the images to obtain the processed burst image  $\bar{b}_1$  and ground truth image  $\bar{y}_{GT}$ . Given this aligned input-ground truth pair, we estimate a pixel-wise color mapping  $C$  between the two images. We assume that the color mapping is linear and model it as a  $3 \times 3$  color correction matrix, which is computed by minimizing a least squares loss. Using the estimated color correction matrix, we can map the network prediction to the same color space as the ground truth and compute pixel-wise error. Our training loss  $\ell(y, y_{GT})$  is thus computed as

$$\ell(y, y_{GT}) = \sum_n m^n \cdot L_1(\hat{y}^n, y_{GT}^n), \quad \hat{y} = C(\phi(y, f_{Pred,GT})) \quad (3)$$

Here,  $\hat{y}$  is the aligned and color mapped network prediction. The summation is over all pixel coordinates  $n$  in the image. The factor  $m^n$  is a binary masking variable used to filter out image regions which are not aligned correctly. It is set to 0 in regions where the error  $R = \|\bar{y}_{GT} - C(\bar{b}_1)\|_2$  after color mapping the processed burst image  $\bar{b}_1$  is greater than a threshold. Note that the images  $\bar{y}_{GT}$  and  $\bar{b}_1$  have lower-resolution compared to the model prediction  $y$ . Thus the error map  $R$  is upsampled to the same resolution as  $y$ , before computing the mask  $m$ .

### 5.3. Training details

We use pre-trained PWC-Net weights for our flow estimator  $F$ . All other modules are initialized using [15]. Our model is first trained using the synthetic data for 100k iterations, and then fine-tuned on the BurstSR dataset for an additional 40k iterations. We use the ADAM [22] optimizer for our training. Data augmentation is performed using random cropping and flipping. Our entire training takes 30 hours on a single Nvidia V100 GPU. All our networks are trained using a burst size of 8.

## 6. Experiments

We perform comprehensive qualitative and quantitative evaluation of our approach. All our experiments are performed for super-resolution by a factor  $s = 4$ . Additional details and results are provided in the suppl. material.

### 6.1. Analysis of our approach

Here, we analyze the impact of different components in the proposed burst SR architecture. We report results on a synthetically generated test set containing 300 bursts, as well as our BurstSR validation dataset. The synthetic test set is generated using the pipeline described in Sec 5.1, using sRGB images from the test set of the Zurich RAW to RGB dataset [17]. We evaluate the networks trained using only the synthetic training data on this set. Since an accurate ground truth HR image is naturally available, the synthetic test set allows us to evaluate the impact of different architectural choices. We also report results on our BurstSR validation set, using the models fine-tuned on BurstSR training set. Since the input burst and HR ground truth in BurstSR are captured using different cameras, there exists spatial and color misalignments between them. We therefore align the network prediction to the ground truth and perform color transformation as described in Sec 5.2. The resulting image is then compared with the ground truth in order to compute performance metrics. We report the standard fidelity based metrics PSNR and SSIM [43], as well as the learned perceptual score LPIPS [49] on both datasets. All metrics are computed in linear sensor space. Note that the images in our BurstSR dataset are generally underexposed, leading to high PSNR scores for all methods. Unless specified, all the methods are evaluated using a burst size of 8.

**Impact of using multiple frames:** Here, we investigate the impact of using multiple frames for SR by comparing our MFSR approach with a single image baseline. We train a SISR network with exactly the same encoder and decoder architecture as employed in our approach. In order to ensure that the SISR performance is not limited by model capacity, we increased the depth of the single image network until its performance saturated. We compare this single image baseline with our multi-frame approach, evaluated using bursts

	Synthetic data			BurstSR		
	PSNR $\uparrow$	LPIPS $\downarrow$	SSIM $\uparrow$	PSNR $\uparrow$	LPIPS $\downarrow$	SSIM $\uparrow$
Single Image	36.42	0.123	0.913	46.41	0.041	0.979
Burst-2	34.90	0.133	0.893	46.10	0.040	0.977
Burst-4	37.18	0.092	0.927	47.06	0.033	0.981
Burst-8	38.61	<b>0.084</b>	0.941	47.52	0.031	0.983
Burst-14	<b>39.09</b>	<b>0.084</b>	<b>0.945</b>	<b>47.76</b>	<b>0.030</b>	<b>0.984</b>

Table 1. Comparison of the baseline SISR network with our multi-frame approach, evaluated using different number of input frames.

	Synthetic data			BurstSR		
	PSNR $\uparrow$	LPIPS $\downarrow$	SSIM $\uparrow$	PSNR $\uparrow$	LPIPS $\downarrow$	SSIM $\uparrow$
Ours	<b>38.61</b>	<b>0.084</b>	<b>0.941</b>	<b>47.52</b>	<b>0.031</b>	<b>0.983</b>
No Alignment	36.66	0.119	0.915	46.50	0.040	0.979
Single Image	36.42	0.123	0.913	46.41	0.041	0.979

Table 2. Comparison of our approach performing explicit alignment with a baseline which does not employ an alignment module.

of different sizes. The result of this comparison is shown in Table 1. Even when using only 4 input frames, our approach significantly outperforms the single image baseline with an improvement of 0.76 dB in PSNR on the synthetic set. Note that although our model is trained using a fixed burst size of 8, it generalizes to bursts with varying input sizes, providing a consistent improvement with increasing burst size. This shows that our approach can effectively utilize the information from multiple frames in order to improve SR performance. When using bursts of size 14, our approach obtains an improvement of 2.67 dB in PSNR on the synthetic set, clearly demonstrating the advantages of using multiple frames for SR.

**Impact of alignment module:** We analyse the impact of the alignment module in our architecture by evaluating a baseline network which does not perform any explicit alignment. We directly concatenate the encoded *base* frame features to all other frames, and pass the resulting feature maps through additional residual blocks, before merging them. The result of this comparison is shown in Table 2. Our approach, performing explicit sub-pixel alignment using a flow estimator, outperforms the baseline **No Alignment** with an improvement of 1.02 dB in PSNR on the BurstSR validation set. Interestingly, the No Alignment network only obtains a slight improvement over the SISR baseline. These results show that accurate alignment of input frames is essential in order to benefit from multiple frames.

**Analysis of fusion architecture:** We compare our proposed attention-based fusion module with 4 different alternatives. i) **MaxPool:** The encoded feature maps are merged by performing element-wise max pooling across the burst. ii) **AvgPool:** The merged feature map is computed as element-wise mean across the burst. iii) **Concatenate:** The encoded feature maps are concatenated along the channel dimension to obtain the merged features. Note that this architecture is constrained to operate on bursts of fixed size. iv) **RecMerge:** The recursive fusion strategy proposed in [6]. Pairs of encoded feature maps are concatenated and

	Synthetic data			BurstSR		
	PSNR $\uparrow$	LPIPS $\downarrow$	SSIM $\uparrow$	PSNR $\uparrow$	LPIPS $\downarrow$	SSIM $\uparrow$
Ours	<b>38.61</b>	<b>0.084</b>	<b>0.941</b>	<b>47.52</b>	<b>0.031</b>	<b>0.983</b>
MaxPool	36.24	0.116	0.912	46.74	0.039	0.980
AvgPool	35.45	0.131	0.902	46.53	0.040	0.979
Concatenate	37.80	0.098	0.928	47.17	0.034	0.981
RecMerge	37.55	0.098	0.927	47.12	0.033	0.981

Table 3. Analysis of different fusion approaches for merging the information from input frames.

	Synthetic data			BurstSR		
	PSNR $\uparrow$	LPIPS $\downarrow$	SSIM $\uparrow$	PSNR $\uparrow$	LPIPS $\downarrow$	SSIM $\uparrow$
Only Feature	37.46	0.101	0.927	47.11	0.034	0.981
Only Residual	38.14	0.093	0.935	47.46	0.031	0.982
Residual+Base	38.41	0.085	0.939	47.46	<b>0.030</b>	<b>0.983</b>
Residual+Base+Flow	<b>38.61</b>	<b>0.084</b>	<b>0.941</b>	<b>47.52</b>	0.031	<b>0.983</b>

Table 4. Impact of different inputs used by the weight predictor.

passed through a small network to merge them. This process is repeated recursively until a single merged feature map is obtained. All four baseline networks employ the same encoder, decoder, and alignment modules as used in our approach to ensure a fair comparison.

The result of this analysis is shown in Table 3. We observe that MaxPool and AvgPool approaches obtain poor results, indicating that simple pooling operations are insufficient to perform effective merging. Both Concatenate and RecMerge achieve better results with PSNR of 37.80 dB and 37.55 dB respectively, on the synthetic set. Our attention-based fusion obtains the best results on both the synthetic set as well as BurstSR, showing that it can effectively merge the information from the input frames.

**Analysis of weight predictor network:** Here, we analyse the impact of different inputs used by our weight predictor network to determine the element-wise fusion weights. We evaluate 4 different versions of the weight predictor, using different sets of inputs, i) **Only Feature:** Only the projected feature map  $\tilde{e}_i^p$  is used. ii) **Only Residual:** Only the feature residual  $r_i = \tilde{e}_i^p - \tilde{e}_1^p$  is used. iii) **Residual+Base:** Both the feature residual  $r_i$  and the base frame features  $\tilde{e}_1^p$  are used. iv) **Residual+Base+Flow:** The feature residual  $r_i$ , base frame features  $\tilde{e}_1^p$ , as well as the flow features  $\hat{f}_i$  are used. The result of this comparison is shown in Table 4. Compared to using only the input feature  $\tilde{e}_i^p$ , using the residuals  $r_i$  instead leads to better performance. Additionally using the base frame features  $\tilde{e}_1^p$  improves the performance further by 0.27 dB in PSNR on the synthetic set. The best results are obtained when using the feature residual  $r_i$ , the base frame features  $\tilde{e}_1^p$ , and the flow features  $\hat{f}_i$  together, showing that they each provide complementary information to the weight predictor.

## 6.2. Comparison with other approaches

In this section, we evaluate our proposed burst super-resolution network on the test set of our BurstSR dataset. We compare our approach with three methods: i) **Single Image** Our SISR baseline network; ii) **DeepJoint+RRDB**

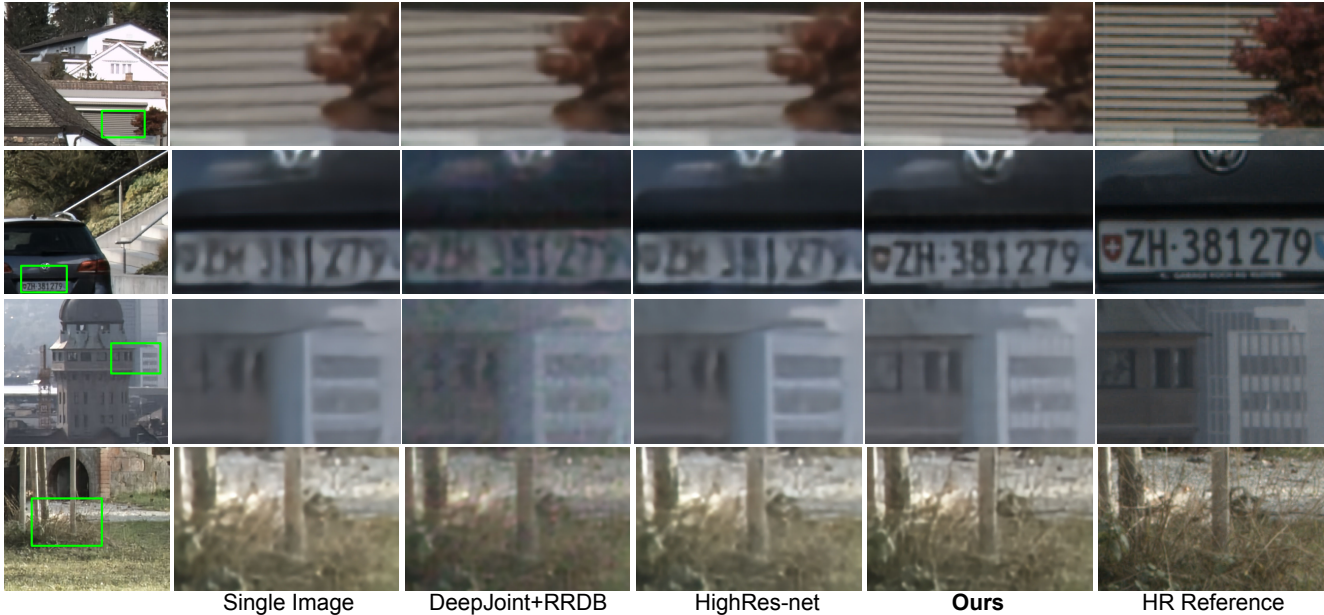


Figure 3. Qualitative comparison of our approach on real world bursts from the BurstSR test set. Our approach can effectively merge information from multiple frames to reconstruct high-frequency image details.

	MOR ↓	%Top ↑	PSNR ↑	LPIPS ↓	SSIM ↑
DeepJoint+RRDB	3.42	8.9	42.13	0.088	0.957
Single Image	2.41	18.3	44.02	0.051	0.972
HighRes-net	2.36	19.2	43.99	0.051	0.972
<b>Ours</b>	<b>1.81</b>	<b>53.6</b>	<b>45.17</b>	<b>0.037</b>	<b>0.978</b>

Table 5. Comparison of our method with existing SR approaches on the BurstSR test set. We report the results of our user study, as well as the standard quality metrics PSNR, LPIPS, and SSIM.

A two stage network which performs single frame demosaicking and denoising using DeepJoint [12] and super-resolves the resulting RGB image using the RRDB [42] network; and iii) **HighRes-net** A recent deep learning based MFSR approach [6] proposed for remote sensing applications. HighRes-net performs implicit registration of the input frames, without using any independent alignment module. Fusion is performed in a recursive manner. We use pre-trained weights for the DeepJoint and RRDB networks. The Single Image baseline, as well as HighRes-net, are trained to perform joint denoising, demosaicking, and SR using the exact training pipeline used by our approach. In order to ensure a fair comparison, we increased the depth of the original HighRes-net network to have the same number of residual blocks as in our approach.

We conducted a user study on Amazon Mechanical Turk to compare the four approaches. We obtain the HR prediction for each of our network on the 20 test images. Next, we extract 15 random  $200 \times 200$  crops from each of our 20 test images. Each of the 300 crops are then resized to  $400 \times 400$  using nearest neighbor interpolation. We show the participants the ground truth HR image, as well as the network predictions. The participants are asked to rank the

predictions from the 4 approaches according to the visual quality w.r.t. the provided DSLR reference image. The network predictions were anonymized and randomized in order to avoid any bias. We obtained 5 independent rankings for each crop. The mean ranking (MOR) over all the crops, as well as the percentage of times a method was ranked first (%Top) are shown in Table 5. Our approach obtains a MOR of 1.81, significantly better than all other approaches. Furthermore, our approach is ranked as the best among all methods 53.6% of the times, more than 2.5 times the second best method. We also report the PSNR, LPIPS, and SSIM scores on the test set, computed as described in Sec. 6.1. A qualitative comparison is also provided in Fig. 3. Our approach obtains the best results in terms of all three metrics, outperforming HighRes-net by 1.18 dB in terms of PSNR.

## 7. Conclusions

We address the problem of real-world multi-frame super-resolution. We introduce a new dataset BurstSR containing RAW burst sequences captured from a handheld camera, and corresponding high-resolution ground truths obtained using a zoom lens. We further propose a multi-frame super-resolution network which can adaptively combine the information from multiple input images using an attention-based fusion. Our approach obtains promising results on real world bursts, outperforming both single frame as well as multi-frame alternatives.

**Acknowledgments:** This work was supported by a Huawei Technologies Oy (Finland) project, the ETH Zürich Fund (OK), an Amazon AWS grant, and an Nvidia hardware grant.



## References

- [1] A. Aitken, C. Ledig, L. Theis, J. Caballero, Zehan Wang, and W. Shi. Checkerboard artifact free sub-pixel convolution: A note on sub-pixel convolution, resize convolution and convolutional resize. *ArXiv*, abs/1707.02937, 2017. 5
- [2] T. Brooks, Ben Mildenhall, Tianfan Xue, Jiawen Chen, Dillon Sharlet, and J. Barron. Unprocessing images for learned raw denoising. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11028–11037, 2019. 5
- [3] Adrian Bulat, J. Yang, and Georgios Tzimiropoulos. To learn image super-resolution, use a GAN to learn how to do image degradation first. In *ECCV*, 2018. 2, 3
- [4] H. Chang, D. Yeung, and Y. Xiong. Super-resolution through neighbor embedding. *Proceedings of the 2004 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2004. CVPR 2004.*, 1:I–I, 2004. 2
- [5] Shengyang Dai, Mei Han, W. Xu, Y. Wu, and Y. Gong. Soft edge smoothness prior for alpha channel super resolution. *2007 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8, 2007. 2
- [6] Michel Deudon, A. Kalaitzis, Israel Goytom, M. R. Arefin, Zhichao Lin, K. Sankaran, Vincent Michalski, S. Kahou, Julien Cornebise, and Yoshua Bengio. Highres-net: Recursive fusion for multi-frame super-resolution of satellite imagery. *ArXiv*, abs/2002.06460, 2020. 1, 2, 7, 8
- [7] C. Dong, Chen Change Loy, Kaiming He, and X. Tang. Learning a deep convolutional network for image super-resolution. In *ECCV*, 2014. 2
- [8] C. Dong, Chen Change Loy, Kaiming He, and X. Tang. Image super-resolution using deep convolutional networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38:295–307, 2016. 1, 2
- [9] Sina Farsiu, Michael Elad, and P. Milanfar. Multiframe demosaicing and super-resolution from undersampled color images. In *IS&T/SPIE Electronic Imaging*, 2004. 2
- [10] M. Fischler and R. Bolles. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Commun. ACM*, 24:381–395, 1981. 5
- [11] W. Freeman, T. Jones, and E. Pasztor. Example-based super-resolution. *IEEE Computer Graphics and Applications*, 22:56–65, 2002. 2
- [12] Michaël Gharbi, Gaurav Chaurasia, Sylvain Paris, and F. Durand. Deep joint demosaicking and denoising. *ACM Transactions on Graphics (TOG)*, 35:1 – 12, 2016. 8
- [13] Daniel Glasner, S. Bagon, and M. Irani. Super-resolution from a single image. *2009 IEEE 12th International Conference on Computer Vision*, pages 349–356, 2009. 2
- [14] R. Hardie, K. Barnard, John G. Bogner, E. Armstrong, and E. Watson. High-resolution image reconstruction from a sequence of rotated and translated frames and its application to an infrared imaging system. *Optical Engineering*, 37:247–260, 1998. 2
- [15] Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 1026–1034, 2015. 6
- [16] H. Hou and H. Andrews. Cubic splines for image interpolation and digital filtering. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 26:508–517, 1978. 2
- [17] A. Ignatov, L. Gool, and R. Timofte. Replacing mobile camera isp with a single deep learning model. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 2275–2285, 2020. 5, 6
- [18] M. Irani and Shmuel Peleg. Improving resolution by image registration. *CVGIP Graph. Model. Image Process.*, 53:231–239, 1991. 2
- [19] Hui Ji and C. Fermüller. Robust wavelet-based super-resolution reconstruction: Theory and algorithm. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31:649–660, 2009. 2
- [20] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *ECCV*, 2016. 1, 2
- [21] Jiwon Kim, J. Lee, and Kyoung Mu Lee. Accurate image super-resolution using very deep convolutional networks. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1646–1654, 2016. 1, 2
- [22] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR*, 2015. 6
- [23] Wei-Sheng Lai, Jia-Bin Huang, N. Ahuja, and Ming-Hsuan Yang. Deep laplacian pyramid networks for fast and accurate super-resolution. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5835–5843, 2017. 1, 2
- [24] C. Ledig, L. Theis, Ferenc Huszár, J. Caballero, Andrew Aitken, Alykhan Tejani, J. Totz, Zehan Wang, and W. Shi. Photo-realistic single image super-resolution using a generative adversarial network. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 105–114, 2017. 1, 2
- [25] Xin Li and M. Orchard. New edge directed interpolation. *Proceedings 2000 International Conference on Image Processing (Cat. No.00CH37101)*, 2:311–314 vol.2, 2000. 2
- [26] Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. Enhanced deep residual networks for single image super-resolution. *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 1132–1140, 2017. 1, 2
- [27] D. Lowe. Object recognition from local scale-invariant features. *Proceedings of the Seventh IEEE International Conference on Computer Vision*, 2:1150–1157 vol.2, 1999. 5
- [28] Xiaoqiang Lu, Haoliang Yuan, Pingkun Yan, Y. Yuan, and X. Li. Geometry constrained sparse coding for single image super-resolution. *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1648–1655, 2012. 2
- [29] Andreas Lugmayr, Martin Danelljan, Luc Van Gool, and Radu Timofte. Srflo: Learning the super-resolution space with normalizing flow. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part V*, volume

- 12350 of *Lecture Notes in Computer Science*, pages 715–732. Springer, 2020. 1, 2
- [30] Andreas Lugmayr, Martin Danelljan, and R. Timofte. Un-supervised learning for real-world super-resolution. *2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW)*, pages 3408–3416, 2019. 2, 3
- [31] Andreas Lugmayr, Martin Danelljan, R. Timofte, Namhyuk Ahn, D. Bai, Jie Cai, Yun Cao, Junyang Chen, Kaihua Cheng, Se-Young Chun, Wei Deng, Mostafa El-Khamy, C. M. Ho, Xiaozhong Ji, A. Kheradmand, Gwantaek Kim, Hanseok Ko, K. Lee, Jungwon Lee, Hao Li, Z. Liu, Zhi-Song Liu, Shuai Liu, Yunhua Lu, Zibo Meng, Pablo Navarrete Michelini, C. Micheloni, Kalpesh P. Prajapati, Haoyu Ren, Yong Hyeok Seo, W. Siu, Kyung-Ah Sohn, Ying Tai, Rao Muhammad Umer, Shuangquan Wang, Huibing Wang, T. Wu, Haoning Wu, Biao Yang, Fuzhi Yang, J. Yoo, Tongtong Zhao, Y. Zhou, Haijie Zhuo, Ziyao Zong, and Xueyi Zou. NTIRE 2020 challenge on real-world image super-resolution: Methods and results. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 2058–2076, 2020. 2, 3
- [32] Andreas Lugmayr, Martin Danelljan, R. Timofte, M. Fritsche, Shuhang Gu, K. Purohit, P. Kandula, Maitreya Suin, A. Rajagopalan, Nam Hyung Joon, Yu Seung Won, Guisik Kim, Dokyeong Kwon, Chih-Chung Hsu, C. Lin, Yuanfei Huang, Xiaopeng Sun, W. Lu, J. Li, X. Gao, and Sefi Bell-Kligler. AIM 2019 challenge on real-world image super-resolution: Methods and results. *2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW)*, pages 3575–3583, 2019. 2, 3
- [33] Andrea Bordone Molini, Diego Valsesia, G. Fracastoro, and E. Magli. DeepSUM: Deep neural network for super-resolution of unregistered multitemporal images. *IEEE Transactions on Geoscience and Remote Sensing*, 58:3644–3656, 2020. 2
- [34] Nhat Nguyen and P. Milanfar. A wavelet-based interpolation-restoration method for superresolution (wavelet superresolution). *Circuits, Systems and Signal Processing*, 19:321–338, 2000. 2
- [35] Shmuel Peleg, D. Keren, and L. Schweitzer. Improving image resolution using subpixel motion. *Pattern Recognit. Lett.*, 5:223–226, 1987. 2
- [36] S. Rhee and M. Kang. Discrete cosine transform based regularized high-resolution image reconstruction algorithm. *Optical Engineering*, 38:1348–1356, 1999. 2
- [37] W. Shi, J. Caballero, Ferenc Huszár, J. Totz, A. Aitken, R. Bishop, D. Rueckert, and Zehan Wang. Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1874–1883, 2016. 4
- [38] Deqing Sun, X. Yang, Ming-Yu Liu, and J. Kautz. PWC-Net: CNNs for optical flow using pyramid, warping, and cost volume. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8934–8943, 2018. 3
- [39] Ying Tai, Jian Yang, and X. Liu. Image super-resolution via deep recursive residual network. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2790–2798, 2017. 1, 2
- [40] R. Tsai and T. Huang. Multiframe image restoration and registration. In *Advances in Computer Vision and Image Processing*, 1984. 2, 5
- [41] E. Ustinova and V. Lempitsky. Deep multi-frame face super-resolution. *arXiv: Computer Vision and Pattern Recognition*, 2017. 2
- [42] Xintao Wang, K. Yu, Shixiang Wu, Jinjin Gu, Yi-Hao Liu, Chao Dong, Chen Change Loy, Yu Qiao, and X. Tang. ESRGAN: Enhanced super-resolution generative adversarial networks. In *ECCV Workshops*, 2018. 1, 2, 8
- [43] Zhou Wang, A. Bovik, H. R. Sheikh, and E. P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13:600–612, 2004. 6
- [44] Valentin Wolf, Andreas Lugmayr, Martin Danelljan, Luc Van Gool, and Radu Timofte. DeFlow: Learning complex image degradations from unpaired data with conditional flows. In *CVPR*, 2021. 3
- [45] B. Wronski, Ignacio Garcia-Dorado, M. Ernst, D. Kelly, Michael Krainin, Chia-Kai Liang, M. Levoy, and P. Milanfar. Handheld multi-frame super-resolution. *ACM Transactions on Graphics (TOG)*, 38:1 – 18, 2019. 1, 2
- [46] Jianchao Yang, Zhaowen Wang, Zhe L. Lin, Xianbiao Shu, and T. Huang. Bilevel sparse coding for coupled feature spaces. *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 2360–2367, 2012. 2
- [47] Jianchao Yang, J. Wright, T. Huang, and Yuliang Ma. Image super-resolution via sparse representation. *IEEE Transactions on Image Processing*, 19:2861–2873, 2010. 2
- [48] X. Yu and F. Porikli. Ultra-resolving face images by discriminative generative networks. In *ECCV*, 2016. 1, 2
- [49] Richard Zhang, Phillip Isola, Alexei A. Efros, E. Shechtman, and O. Wang. The unreasonable effectiveness of deep features as a perceptual metric. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 586–595, 2018. 6
- [50] X. Zhang, Qi feng Chen, R. Ng, and V. Koltun. Zoom to learn, learn to zoom. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3757–3765, 2019. 3
- [51] Yulun Zhang, Yapeng Tian, Y. Kong, B. Zhong, and Yun Fu. Residual dense network for image super-resolution. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2472–2481, 2018. 1, 2
- [52] A. Zomet and Shmuel Peleg. Multi-sensor super-resolution. *Sixth IEEE Workshop on Applications of Computer Vision, 2002. (WACV 2002). Proceedings.*, pages 27–31, 2002. 2