# Fast and Robust Cascade Model for Multiple Degradation Single Image Super-Resolution

Santiago López-Tapia, Nicolás Pérez de la Blanca

*Abstract*—Single Image Super-Resolution (SISR) is one of the low-level computer vision problems that has received increased attention in the last few years. Current approaches are primarily based on harnessing the power of deep learning models and optimization techniques to reverse the degradation model. Owing to its hardness, isotropic blurring or Gaussians with small anisotropic deformations have been mainly considered. Here, we widen this scenario by including large non-Gaussian blurs that arise in real camera movements. Our approach leverages the degradation model and proposes a new formulation of the Convolutional Neural Network (CNN) cascade model, where each network sub-module is constrained to solve a specific degradation: deblurring or upsampling. A new densely connected CNN-architecture is proposed where the output of each sub-module is restricted using some external knowledge to focus it on its specific task. As far we know this use of domain-knowledge to module-level is a novelty in SISR. To fit the finest model, a final sub-module takes care of the residual errors propagated by the previous sub-modules. We check our model with three state of the art (SOTA) datasets in SISR and compare the results with the SOTA models. The results show that our model is the only one able to manage our wider set of deformations. Furthermore, our model overcomes all current SOTA methods for a standard set of deformations. In terms of computational load, our model also improves on the two closest competitors in terms of efficiency. Although the approach is non-blind and requires an estimation of the blur kernel, it shows robustness to blur kernel estimation errors, making it a good alternative to blind models.

*Index Terms*—Single image super-resolution, super-resolution, multiple degradation deconvolution, convolutional neural networks, cascade model.

## I. Introduction

Single Image Super-Resolution (SISR) is a fundamental low-level computer vision that has received considerable attention in recent years [1], [2], [3], [4], [5], [6], [7], [8]. It consists of recovering a high-resolution (HR) image from a given low resolution (LR) image, assuming a degradation model connecting both images. Typically, the general image degradation model assumed in the SISR literature is given by,

$$y = [x \circledast k] \downarrow_s + n, \tag{1}$$

where $x$ is the HR image, $y$ is the LR image, $n$ is the noise, usually additive white Gaussian noise (AWGN), $\downarrow_s$ is the downsampling operator for factor $s$ and $x \circledast k$ represents the convolution of $x$ with the blur kernel $k$.

The current models used to estimate $x$ from $y$ are typically grouped into three categories according to the assumed degradation model and the solver used: interpolation-based (IB), model-based or energy-function optimization (MBO) and learning-based (LB). The IB approach assumes the following simpler model as the degradation model,

$$y = x \downarrow_s + n, \tag{2}$$

and an interpolation technique is proposed as the solver. Clearly, Eq. (2) models the downsampling operation including as part of the noise any other degradation present in the image. Although these methods are the fastest, they are too simple to cope with real blur degradation effects [9], [8], [10].

Traditional MBO approaches [11], [12], [13], [14], [15] define a regularized energy function that is optimized to estimate $x$. The degradation model (1), together with smoothness or edge preservation conditions in solution, are the most common elements that define the energy function [15], [13]. In these models, an optimization is performed for each new LR image, which provides them with high flexibility, under the LR imaging generation conditions, to manage any degradation type. However, closed-form solutions are not feasible and the use of iterative optimization techniques makes MBO approaches computationally expensive. In addition, the hyperparameters must be hand-picked for each degradation type [16].

In recent years, LB approaches have gained significant momentum thanks to the introduction of deep learning (DL). Recently, variable splitting techniques such as half-quadratic splitting (HQS) [17] and alternating direction method of multipliers (ADMM) [18] have been used to incorporate DL in MBO to improve efficiency. In [6], [19], and [20], for instance, a denoising CNN was used to estimate the final HR image. Furthermore, in [8], the authors expand on this approximation by tasking the CNN with upsampling and denoising in a non-blind approach. Clearly, this introduces a connection between MBO and CNN-based approaches because the MBO solution eventually relies on learned modules. However, this connection assumes a new degradation model,

$$y = x \downarrow_s \circledast k^L + n, \tag{3}$$

where blurring $k^L$ is applied on a downsampled image. Until now, no connection has been provided between this new model and the one in Eq. (1).

In contrast to MBO approaches, most LB models do not explicitly use any image degradation model to estimate $x$, assuming that the degradation model is well represented by Eq. (2). A large database of pairs (LR, HR) is used to learn an end-to-end mapping $f(\cdot)$ such that $x = f(y)$ [21], [22], [23],

[24]. The DL models based on Convolutional Neural Networks (CNNs) have shown the highest performance compared to any other approach [2], [4], [5], [7], [25]. The early success of vanilla CNN models in SISR [2], [26] has stimulated the development of deeper designs that have achieved the current state of the art (SOTA) in SISR. As a matter of fact, most of the proposals from the first one, given in [2], have been driven by the incorporation of technical achievements in the CNN-field as for instance, deeper models, and use of residual block to improve accuracy metrics [27], [28], [29], [5], [7], [30], [31], [32]. To improve the performance and processing time, upsampling moves to the last layer of the network [26], [33], [29]. In addition, SRGAN [29] and EDSR [5] improved over those architectures and introduced residual blocks to increase the network depth. DenseSR [30] used residual dense blocks to improve the networks by combining features from different layers. In ESRGAN [32] and RCAN [31] specific blocks to improve performance are proposed, reaching a new SOTA, but at the cost of a significant increase in computing time. Although these models outperform traditional MBO in terms of performance, they lack adaptation to new deformations at test time.

Very recently, CNN-based models (CNN-BMs) assuming Eq. (1) as a degradation model have been proposed. In this scenario, Shocher et al. [34] proposed a technique for the "Zero-Shot" case that exploits the internal recurrence of information inside a single image to train a DL model. This shares the flexibility at test time, but also the high time consumption of MBO approaches. Alternatively, recent works have proposed architectures that encode information about the blur kernel $k$ and introduce it into the network. Riegler et al. [35] used conditioned regression models to encode a Gaussian family of kernels. In SRMD [9], a Principal Component Analysis (PCA) representation of the blur kernel and a stretching strategy to concatenate it with the LR image is proposed. In [10], the latter approach was expanded and improved, with the use of spatial feature transform (SFT) layers [7], to introduce the PCA representation of the blur kernel into the network. By doing so, they are able to implement a very deep residual network with a significant increase in performance over [9]. However, PCA is not a suitable representation to encode families of blurs with different degrees of freedom, in addition to having to be recalculated to cope with new degradation types. Consequently, the CNN models used in [9] and [10] show shortcomings when blur kernels from strong camera movement are present [9].

An alternative to the above CNN-BMs that attempt to solve the SISR problem using a single CNN is to decompose the problem into sub-tasks and use a group of CNNs connected in a cascade fashion to solve it. This "divide and conquer" approach simplifies the function that each sub-module must learn, easing the learning process. If $n = 0$, Eq. (1) can be inverted using an upsampling CNN followed by a deblurring CNN. However, in practice, this approximation does not perform well because of the introduction of strong artifacts during upsampling and the increasing difficulty of deconvolving in HR space [10]. The CNN-BMs based on cascades [10], [8], [36] make the first denoise/deblur and later upsampling

implicitly assuming that Eq. (1) can be approximated by Eq. (3). In [36], an unsupervised network was trained to deblur and denoise the LR image, before using an SR network to estimate the HR. However, as shown in [10] and [8], these cascade approaches underperform compared with other SOTA models. Here, we argue that cascade models can be significantly improved by modifying the connection between the sub-modules and taking advantage of the fact that each sub-module can be specialized in a task by constraining its output using domain knowledge.

In this study, we propose a fast, accurate, and robust CNN cascade-based approach for multiple degradations SISR, CAscade-Deblurring-Upsampling-Fusion (CADUF). Compared with previous cascade-based approaches, our model shows new elements regarding deblur, upsampling, and information processing. In deblurring, the LR image is complemented with the information provided by a deblurred and noisy release of itself obtained from the Weiner filter. We take advantage of the fact that artifacts and noise introduced by the Weiner filter are similar across all possible combinations of blur and images. After extracting and matching features from the two images, we use those features to filter out the noise and artifacts in the Weiner-filtered image. In upsampling, we make use of domain knowledge in the form of "Privileged Information" [37] and constrain the output of upsampling to fulfill an affine projection model. As a result, we obtain an error range that can be reduced by a simple final sub-module. Finally, we improve the communication between sub-modules by using a novel connection approach where each sub-module re-uses the features calculated by all previous sub-modules. Our experiments for non-blind SISR show that CADUF reaches SOTA results being competitive with MBO for unseen blurs but much faster. Furthermore, experiments performed in the blind setting show that our method is also robust to inaccuracies in the blur kernel estimation.

The remainder of the paper is organized as follows. In Section II, we describe the proposed CADUF model and explain each one of its primary components. Section III presents the experimental settings used in this study (dataset, degradations and training hyperparameters) as well as an ablation study of each key component and a comparison with SOTA methods. Finally, Sect. IV concludes the paper.

## II. METHOD

In this section, we first introduce the SR problem settings and propose a new CNN architecture that addresses it by combining multiple CNN sub-modules that solve specific sub-problems. Each of these CNN sub-modules are presented in Sections II-A to II-D. Finally, in Section II-E, we describe in detail the parameters of the proposed architecture.

As previously stated in Section I, we consider that the LR image $y$ is produced from the high-res image $x$ in Eq. (1). Typically, in SISR, $k$ is assumed to be an isotropic Gaussian kernel and $\downarrow_s$ bicubic interpolation [9], [10], but in real applications, different kernel types could arise. In this study, we consider that $k^H$ can also be a more complex blurs, like motion blur. Fig. 1 shows some examples of the complexity of the kernels used during our experiments.
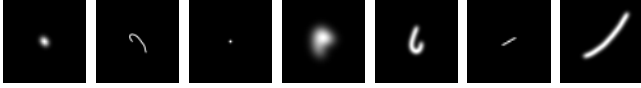
Fig. 1. Examples of the blur kernels $k$ considered in this work. They represent a blur combination of isotropic Gaussian and motion.

We propose a CNN architecture composed of several CNN sub-modules, each of which focuses on solving a portion of the degradation process that produces the LR image. Using this approach has several advantages over using a current SR CNN:

1) By breaking the SR task into sub-tasks each addressed by a sub-module, a strong regularization is imposed over the function space that the model is able to explore. This makes learning easier, leading to improved optima and generalization.
2) Task-specific constraints in each sub-module can be introduced, to further reduce the search space (see Sections II-B and II-C).
3) Specialized architectural elements can be used for each task in each sub-module, improving model performance and efficiency. (see Section II-A).

To do this, we design a cascade based on Eq (3) as an approximation to Eq (1). This allows us to first perform denoising and deblurring before we upsample the LR image. Based on Eq. (3) and assuming that $k^L$ is known, we propose the decomposition of the architecture in the following specialized sub-modules:

1) $E_\theta$: Extraction of motion-corrected features. See Section II-A.
2) $D_\phi$: Deblur and Denoising of the LR image, Section II-B.
3) $U_\psi$: Upsampling, to get an initial SR solution, Section II-C.
4) $F_\omega$: Refinement, improvement of the initial SR solution by cleaning the errors and artifacts caused by the accumulated errors in the cascade. See Section II-D.
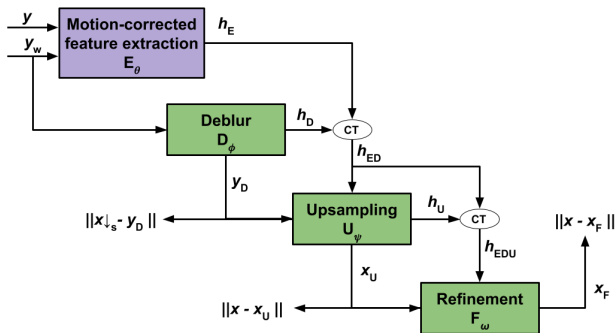


Fig. 2. Proposed architecture for solving the MDSR problem. Each sub-module is a CNN and is specialized in solving a sub-task. Note that "CT" indicates concatenation of the feature maps.

The structure of the proposed architecture is shown in Fig.2. The architecture follows a cascade approach where the features $h$ and the result of each sub-task are needed for the next sub-module in the cascade. Compared to other cascade approaches,

our sub-modules use not only the previously calculated image, but also the features calculated by all previous sub-modules, which allows the propagation and reuse of more information through the cascade. For each spatial location, our model encodes a vector of information, instead of just a value, better representing the distribution of solutions. The architecture is trained in a end-to-end fashion, being its loss function

$$\mathcal{L}_{\text{CADUF}} = \alpha \mathcal{L}_{\text{D}} + \beta \mathcal{L}_{\text{U}} + (1 - \alpha - \beta)\mathcal{L}_{\text{F}} \qquad (4)$$

where $\alpha, \beta > 0$, $\alpha + \beta < 1$ and $\mathcal{L}_{\text{D}}$, $\mathcal{L}_{\text{U}}$ and $\mathcal{L}_{\text{F}}$ are specific losses for sub-modules $D_\phi$, $U_\psi$ and $F_\omega$, respectively. These specific losses will be defined later.

Let us now describe in detail each of the proposed sub-modules of the architecture.

*A. Extraction of motion-corrected features*

Following Eq. (3), the first issue to be addressed is the deconvolution of $y$ by the blur kernel $k^L$ to obtain $x \downarrow_s$. Because the deformations present in the image depend on the combination of the blur kernel $k^L$ and the image, learning a single CNN that can cope with all these variations is a challenging task. For this reason, we propose to start from an initial noisy solution that can be calculated quickly using the Wiener filter:

$$y_w = \mathcal{F}^{-1}\left(\frac{\overline{\mathcal{F}(k^L)}\mathcal{F}(y)}{\overline{\mathcal{F}(k^L)}\mathcal{F}(k^L) + \epsilon}\right), \qquad (5)$$

where $\mathcal{F}$ and $\mathcal{F}^{-1}$ denote the fast Fourier transform (FFT) and inverse FFT, respectively, and $\overline{\mathcal{F}}$ denotes the complex conjugate of $\mathcal{F}$, and $\epsilon$ is the regularization term.

As seen in Section III-E, the use of $y_w$ as the starting point significantly outperforms using only the blur image $y$. However, the Wiener filter increases the noise in the image and introduces artifacts that are difficult to eliminate without the information of the blur image $y$, see Fig 3. To correct the image with a CNN, at each spatial location, we use the features extracted from $y$ and $y_w$. However, because of the blurring processes, $y$ and $y_w$ are not aligned with each other. In this case, we can consider $y_w$ as an anchor that indicates the correct position.
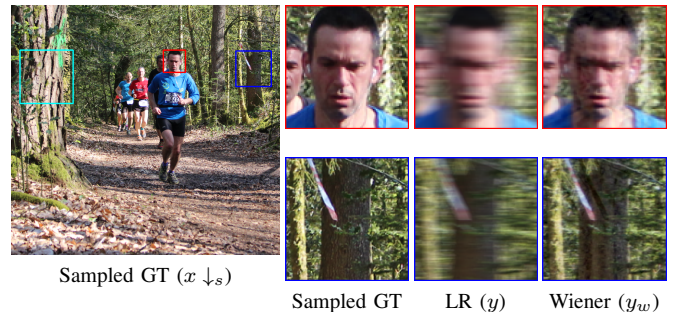


Fig. 3. Illustration of the artifacts originating from using the Wiener filter. Some produced artifacts can be compatible with natural images, being difficult to distinguish from real structures using only the filtered image.

To align the features extracted from both $y$ and $y_w$, we propose the use of deformable convolution. First proposed

in [38] and enhanced in [39] to handle more deformations, deformable convolutions are a modified convolution operation:

$$g^l(m) = \sum_{o=1}^{O} w_o a_{m,o} f^{l-1}(m + o + \Delta m_{m,o}), \qquad (6)$$

where $g^l(m)$ denotes the feature vector in layer $l$ at location $m$, $o$ is the position of the convolutional kernel $w$, and $\Delta m_{m,o}$ and $a_{m,o}$ are offsets and scalar modulation, respectively. These parameters are calculated for each $m$ location using another convolutional layer. Owing to this operation, our model can align the features of $y$ and $y_w$ at each spatial location, making further downstream calculations much easier.
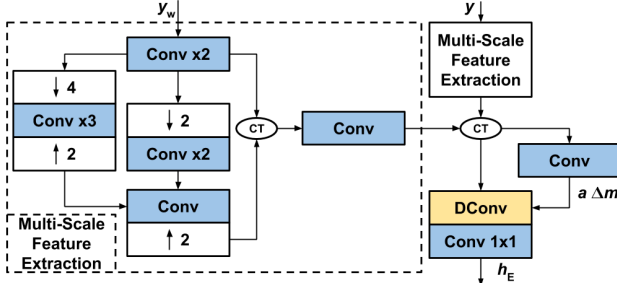


Fig. 4. Motion-corrected feature extraction sub-module. Conv is a convolution, and DConv is a deformable convolution. Each convolution is followed by a leaky-ReLU activation with 0.2 of negative slope and has a $3 \times 3$ kernel with the exception of the first and last convolutions that use a $5 \times 5$ and $1 \times 1$, respectively. Convolutions for scale 1 use 64 features while convolutions on scales /2 and /4 use 32.

Nevertheless, the calculation of parameters $a$ and $\Delta m$ is not an easy task: Movement of the camera can be quite strong, thus matching locations between $y$ and $y_w$ can be too far away. In these cases, a large receptive field is required. To obtain this, we extract features at scales 1, /2, and /4 using the architecture shown in Fig. 4. By doing so, we increase the receptive field to 37 pixels without the need to use a far deeper architecture.

This sub-module, $E_\theta$, produces the features $h_E$ that will be used by the remaining network. Although it could move the feature vectors anywhere, by propagating the loss of the next sub-modules in the cascade, it would have to learn to put then in the correct place, producing motion-corrected features.

### B. LR-Anchor image estimation

According to Eq. (3), we start mapping the LR image $y$ to an LR image $y_D^*$ without blur or noise.

To define $y_D^*$, one approach can be defining it as $x \downarrow_s$. However, the difficulty of the SR task primarily depends on the blur kernel used before downsampling. Those cases with both low and high Gaussian blur led to significantly worse reconstruction. Large kernels in HR produce highly blurry LR images while small kernels in HR do not have a smooth downsampling form in LR, which introduces aliasing artifacts in the LR image. Aliased LR images can be very difficult to super-resolve because artifacts must be distinguished from real high-frequency patterns in the image.

Following the above arguments, we define $y_D^*$ as the LR image obtained by downsampling $x$ with a downsampling operator $\downarrow_s^D$ that it is easier to invert by the upsampling network that we use. That is:

$$y_D^* = x \downarrow_s^D = [x \circledast k^D] \downarrow_s, \qquad (7)$$

where $\downarrow_s$ is the bicubic downsampling of factor $s$, and $k^D$ is an isotropic Gaussian with $\sigma$ 0.8 and 1.8 for scaling factors 2 and 4, respectively. These values have been found experimentally from within the range [0.2, 4.0]. Notice that the selection of these values can be seen as a form of "Privileged Information" [37]. The images $\{y_D^*\}$ represent additional information available only during training and are used to estimate the optimal LR image. In this regard, $y_D^*$ acts as an anchor that guides and regularizes the model training.

The features for the deblurring task are computed using a CNN-regression model that maps its input $\{h_E, y_w\}$, the motion-corrected features, and the low-resolution Wiener deconvolved image, into $y_D^*$, the downsampled version of $x$. Let us denote by $D_\phi(h_E, y_w)$ the mapping defined by the CNN model. Then, $D_{\hat\phi}(h_E, y_w) : \{y_w\} \rightarrow \{y_D\}$, where $\hat\phi$ represents the weights learned after training the model. The parameters can be optimized introducing the following loss $\mathcal{L}_D$

$$\mathcal{L}_D = \frac{1}{N} \sum_{i=1}^{N} \mathcal{L}(D_\phi(h_E, y_w), y_D^*), \qquad (8)$$

where $\mathcal{L}(\cdot)$ is any loss function that estimates the differences between two images, like mean squared error (MSE). It can be observed $D_\phi$ computes a deblurring and denoising mapping for the degradation model defined in Eq. (3) because the low-resolution image $y$ is mapped into an image without blur or noise, $y_D$, that is easier to super-resolve.

### C. Initial high-res image estimation

After deconvolution, we are left with upsampling. This is equivalent to the problem defined in Eq. 2 when $n = 0$ and $\downarrow_s^D = \downarrow_s$. Let us define the minimization of the upsampling process as follows:

$$\hat\psi = \arg\min_\psi \frac{1}{N} \sum_{i=1}^{N} \mathcal{L}(P_\psi(z), x), \qquad (9)$$

where $P_\psi$ is a CNN with parameters $\psi$, and $z$ is the input that will be defined later. It should be noted that this approach is similar to the one used in other cascade methods [10], [8]. If the error introduced by the previous steps in the cascade is low, it is expected that it will not significantly affect this step. In this case, the simplest solution, therefore, would be to use as an upsampling sub-module a classic CNN for SISR such as SRResNet [29] to minimize Eq. 9. However, as shown by our experiments in Section III-E, this is not the best solution in practice. By minimizing Eq. 9, it has to solve two distinctive tasks: Inversing the operator $\downarrow_s^D$ and removing the approximation error. The learning of the parameters is more difficult: Small differences can be due to approximation errors or different structures in the HR space. If not constrained, small changes in the LR image can lead to very different SR

images, leading to an ill-conditioned model. This could result in the appearance of strong artifacts.

To use Eq. 9 to learn only the inverse of $\downarrow_s^D$, we constrained the predicted image $x_U$ to those solutions that satisfy,

$$x_U \downarrow_s^D = y_D, \tag{10}$$

through the use of the projection introduced in [40]. Let us define $A$ as a degradation operator such that $Ax = x \downarrow_s^D$ and $A^+$ as the Moore-Penrose pseudoinverse [41] of $A$ (refer to Section III-B to see the estimation of $A^+$). In our case, $A$ is an affine projection implemented as a stride convolution, and $A^+$ is implemented using a convolutional transpose layer. For a given $P_\psi(z)$, a new transformation $U_\psi(z)$ can be defined as,

$$U_\psi(z) = (I - A^+ A)P_\psi(z) + A^+ y_D, \tag{11}$$

such that if no noise is present because $AA^+A = A$ and $A$ is a full row rank matrix, that is, $AA^+ = I$, then,

$$AU_\psi(z) = A(I - A^+ A)P_\psi(z) + AA^+ y_D = 0 + y_D. \tag{12}$$

Therefore, if $x_U = U_\psi(z)$, Eq. 10 is always satisfied.

Owing to the loss of information caused by approximation errors, the information of $y_D$ contained in $h_D$ is not sufficient to properly invert $\downarrow_s^D$. Our model needs information of the original $y$, which could not be preserved in $h_D$. We add this information by concatenating to $h_D$ the motion-corrected features $h_E$, obtaining $h_{ED}$. Note that this would not be possible without motion-correction because the features $h_D$ will not align with features extracted from the blurred image $y$. Therefore, we define the input as $z = \{h_{ED}, y_D\}$. Using these inputs and the affine projection, we obtain a specialized SR sub-module by replacing $P_\psi$ for $U_\psi$ and introduced the following loss to optimized its parameters:

$$\mathcal{L}_U = \frac{1}{N} \sum_{i=1}^{N} \mathcal{L}(U_\psi(h_{ED}, y_D), x). \tag{13}$$

### D. High-res image refinement

Compared to previous SISR cascade models, ours adds an improvement in the form of a last fusion block fed with the characteristics of all the blocks in the cascade (see Fig. 2) focused on removing the accumulated error. This last block is the result of distinguishing two tasks in the previous step: Inverting the downsampling operator $\downarrow_s^D$, in addition to removing the accumulated error $e$ introduced by the previous cascade sub-modules,

$$x = x_U + e. \tag{14}$$

Note that $e$ can be seen as being structured noise dependent not only on $y$, but also on previous sub-modules $E_\theta$, $D_\phi$, and $U_\psi$. Because $U_\psi$ is constrained, $e$ will not introduce strong artifacts when the error introduced by $D_\phi$ is small (as shown in Section III-E). Therefore, we argue that a last refinement sub-module $F_\omega$, when fed with the features $h_{EDU}$ calculated by each one of the specialized sub-modules, $E_\theta$, $D_\phi$, and $U_\psi$, can learn to calculate a residual that, when added to $x_U$, compensates for the accumulated errors $e$ in the cascade.

We implemented $F_\omega$ using a CNN regression model. The learning of the parameters is performed using the following loss function:

$$\mathcal{L}_F = \frac{1}{N} \sum_{i=1}^{N} \mathcal{L}(F_\omega(h_{EDU}, x_U), x). \tag{15}$$

### E. $D_\phi$, $U_\psi$ and $F_\omega$ architecture

The architecture that implements our sub-modules $D_\phi$, $U_\psi$, and $F_\omega$ is shown in Fig. 5. As seen, it is a modification of the residual architecture SRResNet proposed in [42]. We remove the batch normalization layers as suggested in [5] and replace the ReLU activation with leaky ReLU with 0.2 negative slope. To further increase the receptive field, the first convolution in each residual block is a dilated convolution [43] with the dilatation parameter set to 2. The architecture is fed with features $h$ calculated by the previous sub-modules. The first two convolutional layers transform the input features. Then, $Nb$ residual blocks are used to produce the new set of features $\hat{h}$. To calculate the output image $\hat{z}$, we use $\hat{h}$ together with a dynamic filter network (DFN) [44]. The DFN predicts $c$, which is a collection of filters, one for each spatial location in the output image. Using the approximation proposed in [45], $\hat{z}$ is calculated as,

$$\hat{z}_{i,j} = r_{i,j} + \sum_{l=-p}^{p} \sum_{m=-p}^{p} c_{i,j,l,m} * z_{i+l,j+m}, \tag{16}$$

where $z$ is the output image calculated by the previous sub-module, $i,j$ is a spatial location, and $r$ is a residual image calculated by a CNN using $\hat{h}$. We use a value of $p = 2$. By using this approximation, our model can filter spatially-variant artifacts present in the outputs of previous sub-modules. Then, the filtered outputs can be used to perform residual learning [46], [47] more effectively. These artifacts are caused by approximation errors in the image formation model and the calculation. Notice that a normal CNN can also learn to remove these artifacts, but it would require far more parameters because CNN filters are spatially-invariant.

In each convolution, except in the input convolutions of each sub-module and the two output convolutions, kernels of size $3 \times 3$ and 64 filters are used. The input convolutions use kernels of size $1 \times 1$ to fuse the features of the previous sub-modules. In our experiments, the number of residual blocks $Nb$ was set to eight for $D_\phi$, 16 for $U_\psi$, and $F_\omega$ uses only four. In the case of $U_\psi$ and $F_\omega$, both $c$ and $r$ are upsampled by scaling factor $s$ using sub-pixel convolutions [33]. The sub-pixel convolution is located after the last convolution for $c$ and before in the case of $r$. The first two convolutions of $D_\phi$ are skipped because $E_\theta$ already provides a suitable set of features.

## III. EXPERIMENTS

### A. $k^{*L}$ estimation

To train our models, for each $k$ in Eq. (1), we need to have $k^{*L}$ for Eq. (3) such as $k^{*L} = \arg\min_{k^L} ||[x \circledast k] \downarrow_s - x \downarrow_s \circledast k^L||_F^2$, that is, $k^{*L}$ minimizes the difference between both degradation models. To calculate $k^{*L}$, we use a neural network
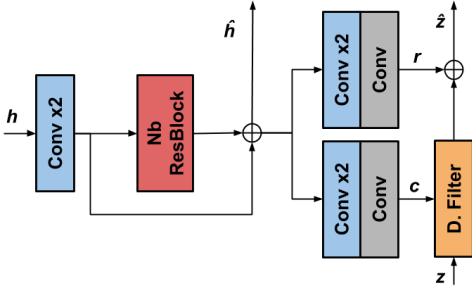
Fig. 5. Architecture used for sub-modules $D_\phi$, $U_\psi$, and $F_\omega$. Conv indicates a convolutional layer and D. Filter the correlation of kernels $c$ with an input image $z$. We use Nb residual blocks [42], each composed of two convolutions after removing the batch normalization layers, as suggested in [5]. Each convolution uses a $3 \times 3$ kernel size, with the exception of the first layer that uses $1 \times 1$ kernels to fuse the previous sub-module features. With the exception of the output convolutions, all are followed by a leaky ReLU activation with a negative slope of 0.2. We skip the first two convolutions for $D_\phi$ because $E_\theta$ already provides a suitable set of features. In the case of $U_\psi$ and $F_\omega$, sub-pixel convolutions [33] of scaling factor $s$ are used to produce $c$ and $r$ with the correct output size (see the text for further details).

with one hidden layer with 2048 neurons. This network $L_\zeta$ with parameters $\zeta$ is trained by minimizing $||x \circledast k| \downarrow_s - x \downarrow_s \circledast L_\zeta(k)||_F^2$. We use the Adam optimizer for 60 epochs with lr $= 10^{-4}$ and weight decay set to $10^{-4}$.

### B. $A^+$ calculation

To use the affine projection proposed in Section II-C, we implement the $A^+$ operator using a CNN with two convolutional layers with 32 filters followed by leaky ReLU ($\alpha = 0.2$) activation and a final sub-pixel convolution. The sizes of the kernels are $7 \times 7$, $5 \times 5$, and $3 \times 3$, respectively. Following the approach in [40], we train this network by minimizing the following loss function with a stochastic gradient descent, that is,

$$\hat{\mu} = \arg\min_\mu \mathbb{E}_x \|Ax - AA_\mu^+(Ax)\|_2^2$$
$$+ \mathbb{E}_y \|A_\mu^+(y) - A_\mu^+(AA_\mu^+(y))\|_2^2, \quad (17)$$

where $A$ is the degradation operator and $A_\mu^+$ is a CNN of the parameters $\mu$. We stop the optimization when the loss value is less than $10^{-7}$. During the training of our main network, we keep $\mu$ fixed.

### C. Datasets

The training dataset is formed by images from the DIV2K [50] and Flickr2K [51] datasets. The ground-truth of our training consists of 3450 high-quality 2K images. During training, random patches cropped from the images of the dataset were extracted. The size of each HR patch is $s48 \times s48$, where $s$ is the scaling factor. The LR images were synthesized according to Eq. (1). For testing, we used images from standard SR test image datasets BSD100 [48], Urban100 [3], and Manga109 [49]. To ensure that a wide variety of degradations are present in the test, we generated three LR images for each of the HR images.

The blur kernels used to generate the LR images cause a uniform blur and are a combination of an isotropic Gaussian

blur kernel and a motion blur kernel. We simulate two different motion blur kernels representing simple and complex camera movements. In the simple case, only simple linear movement is considered. Meanwhile, the motion blur kernels for the complex case simulate more complex camera movements with curved trajectories. They are generated using the method in [52] with $T = 0.8$ and anxiety $10^R/1000$, where $R$ is a random number from a uniform $[0, 1]$ distribution. Furthermore, in this scenario, AWGN of standard deviation $\sigma = 0.01$ is added to the LR image. The combination of Gaussian blur and the motion blur kernels of each scenario produces two sets of blur kernels that we use in our experiments. Fig.1 in Columns (1,3,6) and (2,4,5,7) show blur kernel examples of the simple and complex case respectively. We refer to these two sets of kernels, GaussianSM and GaussianCM, respectively.

In the case of GaussianSM, the $\sigma$ of the Gaussian blur is sampled randomly during training from the range [0.2, 2.0] and [0.2, 4.0] for scaling factors 2 and 4, respectively. The motion blur kernels of GaussianSM are sampled with angles in the $[0, 180)$ range and length in the range $[1, 9]$ and $[1, 15]$ for scaling factors 2 and 4, respectively. GaussianCM consists of 640 blur kernels with sizes between $11 \times 11$ and $45 \times 45$ pixels that are combined with a Gaussian blur with $\sigma$ in the range [0.2, 1.0] and [0.2, 2.0] for factors 2 and 4, respectively. We used 540 of these kernels for training and 100 for testing.

### D. Model training

We train our models using the Adam optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.999$, weight decay of $10^{-4}$, and a batch size of 64. For each epoch, we sample 3000 batches. We define L as the Charbonnier loss [53], as we found it to be more stable than MSE. The Charbonnier loss between two images $u$ and $v$ is $\gamma(u, v) = \sum_i \sum_j \sqrt{(u_{i,j} - v_{i,j})^2 + \varepsilon^2}$, where $\varepsilon = 10^{-3}$. We apply data augmentation by randomly performing horizontal and vertical flips, 90-degree rotations, and random scaling with a factor between [0.5, 1.25]. The parameter of the Wiener filter is set as the minimum between 0.001 and the estimated standard deviation of the noise in the images. All of our models use the RGB images as input and do not require any additional preprocessing step. The implementation was performed using Pytorch [54][1].

Our training method consists of two distinctive phases: A) An initialization phase where we train our model using much higher $\alpha = 0.6$ and $\beta = 0.3$, forcing the model to first converge to good solutions for the intermediary problems of an optimal LR image and initial high-res image estimation. Although it is possible to train the proposed model without this initial step, we have found that this initialization helps stabilize the model and makes it converge earlier. We trained for 20 epochs with lr $= 10^{-4}$.

B) Main training phase. We set $\alpha = 0.1$ and $\beta = 0.1$ and train our model for 120 epochs. The learning rate was set to $10^{-4}$ for the first 90 epochs and then set to $10^{-5}$ for the remaining 30 epochs.

---

[1]The code and trained models shown in this work will be made available upon acceptance of the paper at https://github.com/vipgugr/CADUF.

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|
| Model | $P_\psi(y)$ | $P_\psi(y,\mathrm{PCA}(k))$ | $P_\psi(y_w)$ | $P_\psi D_\phi(y_w)$-no-PI | $P_\psi D_\phi(y_w)$ | $F_\omega U_\psi D_\phi(y_w)$ | $F_\omega U_\psi D_\phi(y,y_w)$ | CNN-Cascade | CADUF |
| PSNR/SSIM | 26.63-0.773 | 27.06-0.780 | 27.09-0.785 | 27.24-0.789 | 27.37-0.792 | 27.51-0.795 | <u>27.53-0.796</u> | 27.25-0.788 | **27.90-0.806** |

TABLE I

KEY COMPONENT ANALYSIS OF THE PROPOSED MODEL FOR SCALE FACTOR 4 ON THE GAUSSIANSM TESTING DATASET (IMAGES FROM BSD100 [48], URBAN100 [3] AND MANGA109 [49] DATASETS). BEST AND SECOND-BEST RESULTS APPEAR IN BOLD AND UNDERLINED, RESPECTIVELY.

*E. Ablation study*

In this section, we conducted experiments to determine the contribution of each component of the proposed CADUF. All models have 28 residual blocks and are trained as described in Section III-D. The experiments were conducted for scale factor 4 and using GaussianSM training and testing sets (see Section III-C),with the difference that the blurs used were sampled with $\sigma$ in the range of $[0.2, 3.0)$ and length of $[1, 11]$. This was done to ease the training on the simpler models because using strong blurs makes it unstable. Table I contains the results for this study in terms of the PSNR and SSIM.

The base model $P_\psi(y)$ is a single CNN model that uses the architecture described in Section II-E (see Fig 5), where $h$ is the features extracted from $y$ with two convolutional layers and $z = y$. This model shows the performance that can be expected from a classic SR CNN model. We improve the model performance by progressively adding each main component of the CADUF.

The first improvement of 0.46dB-0.12(SSIM) is achieved by introducing the $k^L$ blur information into the network using the Weiner-filtered image $y_w$ instead of $y$ ($P_\psi(y_w)$) (Columns 2 and 4 in Table.I). Facilitating the blur information eases the network task because the blur no longer needs to be identified. Alternatively, Zhang et al. [9] proposed using PCA to encode the information of the blur kernel $k$ and use it alongside the blur image $y$. We also test this approximation, making the model $P_\psi(y,\mathrm{PCA}(k))$ (Column 3 in Table.I). Despite the noise and artifacts introduced by the Weiner filter and the approximation error of using Eq. (3), $P_\psi(y_w)$ slightly outperforms $P_\psi(y,\mathrm{PCA}(k))$. Moreover, the use of PCA requires its training with the possible blurs $k$, which harms the generalization capabilities of the model to unseen blurs.

The second improvement is the decomposition of $P_\psi(y_w)$ into $P_\psi D_\phi(y_w)$ (Columns 4 and 6 in Table.I). By specializing part of the model in the denoising and deconvolution of the LR image, we can see an increase in performance of 0.28dB-0.007(SSIM) despite the use of the same number of residual blocks (a total of 28 residual blocks, eight residual blocks for $D_\phi$ and 20 for P). We also tested the contribution of using privileged information (PI) [37] during training. By training $D_\phi$ to map the input to $x \downarrow_s$ instead of $x \downarrow_s^D$, we obtain a new model that does not make use of PI. We call this model $P_\psi D_\phi(y_w)$-no-PI. As seen in Columns 5 and 6 of Table I, the use of PI is a key element of the proposed CADUF because omitting it significantly deteriorates the performance (0.13dB-0.003(SSIM)).

Next, we use the approximation proposed in Section II-C to separate $P_\psi$ into $F_\omega U_\psi$, obtaining the model $F_\omega U_\psi D_\phi(y_w)$ (Column 7-Table I). Using the same number of residual blocks (28), $F_\omega U_\psi D_\phi(y_w)$ outperforms $P_\psi D_\phi(y_w)$ (Column 6-Table I) by 0.14dB-0.003(SSIM). This is in line with our hypothesis that the decomposition of the task of $P_\psi$ into two sub-tasks, upsampling, and error correction, simplifies the optimization and allows for better minima. It can also be observed that, owing to the introduction of prior information and the use of constraints, the concatenation of $D_\phi$ and $U_\psi$ does not output complex artifacts that are difficult to filter out by $F_\omega$. As shown in [55], CNN architectures introduce deep priors, which in some applications explain a large portion of the model efficiency. However, in our case, the ablation study shows that it is the specific training, with the introduction of $F_\omega U_\psi$ and privileged information, that is the primary factor responsible for the model efficiency.

We also tested the proposed mechanism to connect all the sub-modules of the CADUF. To do that, we develop a new model by feeding each sub-module only with the image $z$ calculated by the previous one (not the features $h$) and omitting the reuse of $z$ with dynamic filtering. By doing this, we obtain a classic CNN-cascade model similar to those used in the literature [36], [10], a concatenation of a deblur and upsampling CNN. This model, which we call CNN-Cascade (Column 9 in Table.I), obtains results far worse than $F_\omega U_\psi D_\phi(y_w)$ and even $P_\psi D_\phi(y_w)$ (0.26dB-0.007(SSIM) and 0.12dB-0.004(SSIM), respectively). As we argue in Section II, the difference in performance is because the use of the features of previous sub-modules allows for better information propagation through the cascade.

Finally, we evaluate the contribution of our proposed motion-corrected feature extraction sub-module, $E_\theta$, appending the information extracted from $y$ and $y_w$ together. Overall, it constitutes the full CADUF model (Column 10 in Table.I). Compared to the concatenation of $y$ and $y_w$ ($F_\omega U_\psi D_\phi(y,y_w)$) (Column 8 in Table.I), the use of $E_\theta$ significantly increases the performance (0.37dB-0.010(SSIM)) with little overhead computation.

*F. SOTA Comparison*

In this section, we compare our CADUF model against the SOTA methods in SISR dealing with multiple degradations. RCAN [56] is a very deep CNN developed for SISR that does not use any information about the degradation but achieves a very high score. In contrast, SRMDNF [9] and SFTMD [10] incorporate information of the kernel $k$ into the network by encoding it using PCA, but obtaining lower performance than RCAN [56]. We also show the results of combining SOTA deconvolution methods with an SISR CNN in a cascade. Specifically, we compared two combinations, both using deconvolution followed by upsampling: IRCNN [6]+RCAN

| Scale | Model | PSNR-SSIM | | | | | | MACs |
|---|---|---|---|---|---|---|---|---|
| | | GaussianSM | | | GaussianCM | | | |
| | | BSD100 [48] | Urban100 [3] | Manga109 [49] | BSD100 [48] | Urban100 [3] | Manga109 [49] | |
| ×2 | Bicubic | 26.21-0.692 | 23.58-0.685 | 25.54-0.809 | 23.24-0.543 | 20.68-0.529 | 21.28-0.655 | $2.02 * 10^7$ |
| | ZSSR [34] | 26.55-0.717 | 23.81-0.711 | 25.76-0.810 | 24.11-0.556 | 21.79/0.591 | 21.92-0.664 | $4.70 * 10^{13}$ |
| | RCAN [56]* | 30.85-0.874 | 29.22-0.882 | 34.27-0.956 | 25.13-0.640 | 22.61-0.639 | 24.40-0.775 | $2.46 * 10^{12}$ |
| | IRCNN [6]+RCAN [56] | 27.08-0.759 | 24.66-0.751 | 27.92-0.879 | 26.23-0.712 | 24.30-0.732 | 27.01-0.846 | $2.70 * 10^{12}$ |
| | DeblurGAN+RCAN [56] | 26.21-0.702 | 23.57-0.694 | 25.54-0.815 | 23.99-0.573 | 21.55-0.568 | 22.86-0.706 | $2.47 * 10^{12}$ |
| | SRMDNF [9]* | 30.83-0.873 | 29.42-0.884 | 34.94-0.957 | 25.99-0.695 | 24.01-0.712 | 26.36-0.825 | $\mathbf{2.60 * 10^{11}}$ |
| | SFTMD [10]* | 31.57-0.889 | 30.62-0.907 | 36.79-0.969 | 27.26-0.742 | 25.66-0.774 | 28.97-0.880 | $9.71 * 10^{11}$ |
| | DPSR [8] | 28.37-0.793 | 27.11-0.815 | 33.19-0.925 | 27.23-0.740 | 25.89-0.790 | 30.11-0.895 | $4.86 * 10^{12}$ |
| | CADUF | **31.92-0.893** | **31.11-0.914** | **37.70-0.972** | **27.92-0.761** | **26.64-0.805** | **30.73-0.906** | $6.56 * 10^{11}$ |
| ×4 | Bicubic | 24.19-0.582 | 21.51-0.564 | 22.52-0.698 | 23.05-0.533 | 20.40-0.511 | 20.99-0.642 | $1.68 * 10^7$ |
| | ZSSR [34] | 24.76-0.589 | 21.83-0.573 | 22.88-0.709 | 23.32-0.545 | 20.90-0.520 | 21.45-0.652 | $2.14 * 10^{13}$ |
| | RCAN [56]* | 26.92-0.713 | 24.86-0.736 | 28.80-0.877 | 24.80-0.610 | 22.35-0.612 | 24.13-0.756 | $6.38 * 10^{12}$ |
| | IRCNN [6]+RCAN [56] | 24.71-0.596 | 22.18-0.587 | 23.60-0.714 | 24.19-0.567 | 21.87-0.565 | 23.21-0.691 | $6.98 * 10^{12}$ |
| | DeblurGAN+RCAN [56] | 23.48-0.550 | 21.34-0.538 | 22.31-0.663 | 22.11-0.498 | 20.59-0.494 | 21.35-0.618 | $6.40 * 10^{11}$ |
| | SRMDNF [9]* | 26.82-0.708 | 24.76-0.729 | 28.53-0.870 | 25.31-0.632 | 23.10-0.649 | 25.45-0.793 | $\mathbf{6.71 * 10^{10}}$ |
| | SFTMD [10]* | 27.12-0.721 | 25.25-0.752 | 29.42-0.889 | 25.67-0.647 | 23.62-0.674 | 26.46-0.821 | $2.82 * 10^{11}$ |
| | DPSR [8] | 25.06-0.641 | 23.47-0.686 | 27.35-0.842 | 24.79-0.629 | 22.51-0.656 | 24.95-0.804 | $1.33 * 10^{12}$ |
| | CADUF | **27.41-0.729** | **25.64-0.768** | **30.05-0.900** | **25.73-0.652** | **23.82-0.689** | **26.78-0.835** | $2.03 * 10^{11}$ |

TABLE II

COMPARISON WITH NON-BLIND SOTA ON THE BSD100 [48], URBAN100 [3] AND MANGA109 [49] DATASETS FOR THE BLUR KERNELS IN GAUSSIANSM AND GAUSSIANCM DATASETS. TO ENSURE FAIRNESS IN THE COMPARISON, THE METHODS INDICATED WITH "*" WERE RETRAINED. THE AVERAGE NUMBER OF MACS PER IMAGE OF EACH METHOD IS ALSO REPORTED. BEST AND SECOND-BEST RESULTS ARE BOLD AND UNDERLINED, RESPECTIVELY.

[56] and DeblurGAN+RCAN [56]. Finally, we also compare the MBO approaches ZSSR [34] and DPSR [8].

We test all models using the test datasets GaussianSM and GaussianCM described in Section III-C. To ensure fairness in the comparison, we not only train our model for each setting using the appropriate training set, but also do the same for competing models that require it. Therefore, we trained RCAN [56], SRMDNF [9], and SFTMD [10]. The training and evaluation of the competing methods was performed using the code provided by the original authors.

Table II shows the results of our CADUF model and competing models in the non-blind scenario for GaussianSM and GaussianCM kernels datasets and scaling factors 2 (x2) and 4 (x4). We evaluate the performance of the methods using the PSNR and SSIM metrics. Furthermore, we also include a study of the time complexity of each method using the average number of multiplication and additions (MACs) per image as a metric. The smaller the number of MACs, the faster the algorithm performance.

As seen in Table II, compared to other LB methods, CADUF performs significantly better in all datasets and scenarios. In GaussianSM, CADUF obtains an average improvement on the datasets of 2.13dB-0.02(SSIM) in x2 and 0.84dB-0.02(SSIM) in x4 over RCAN [56] and 0.58dB-0.01(SSIM) in x2 and 0.44dB-0.01(SSIM) in x4 over SFTMD [10], being significantly faster. The performance difference increases for the GaussianCM dataset. In this case, the figures show improvements of 4.38dB-0.14(SSIM) in x2 and 1.68dB-0.06(SSIM) in x4 over RCAN [56] and 1.14dB-0.03(SSIM) in x2 and 0.19dB-0.01(SSIM) in x4 over SFTMD [10].

The difference between CADUF and the other best performing methods in GaussianSM is more apparent if we perform a closer inspection of the produced images in Fig. 6. As shown, when uniform motion blur is present in the image, CADUF is

| Scale | Model | PSNR-SSIM | | |
|---|---|---|---|---|
| | | BSD100 [48] | Urban100 [3] | Manga109 [49] |
| ×2 | Bicubic | 23.24-0.543 | 20.68-0.529 | 21.28-0.655 |
| | RCAN [56]* | 25.13-0.640 | 22.61-0.639 | 24.40-0.775 |
| | DeblurGAN+RCAN [56] | 23.99-0.573 | 21.55-0.568 | 22.86-0.706 |
| | IKC [10]+SFTMD [10]* | 18.27-0.451 | 16.29-0.456 | 17.50/0.585 |
| | DPSR [8] | 23.56-0.678 | 20.32-0.652 | 23.57-0.825 |
| | Our Model | **25.65-0.698** | **23.27-0.674** | **25.06-0.847** |

TABLE III

COMPARISON WITH BLIND SOTA ON THE BSD100 [48], URBAN100 [3] AND MANGA109 [49] DATASETS FOR THE BLUR KERNELS IN THE GAUSSIANCM DATASET. TO ENSURE FAIRNESS IN THE COMPARISON, THE METHODS INDICATED WITH "*" WERE RETRAINED. BEST AND SECOND-BEST RESULTS ARE BOLD AND UNDERLINED, RESPECTIVELY.

able to recover fine details in the image as well as the other best performing method, SFTMD [10], but without introducing ghosting and other artifacts. Despite being a blind method, RCAN [9] is able to compete with the remaining methods in GaussianSM, obtaining similar results to other non-blind methods such as SRMDNF [9]. However, its performance drops significantly for the more complex blurs in GaussianCM, showing that even a CNN as complex as RCAN [9] is not able to cope with the degradation variance of a more complex realistic scenario without the introduction of specialized elements. In this case, although SRMDNF [9] and SFTMD [10] do not suffer as much as RCAN [9], they show a more significant drop in performance than CADUF and DPSR [8]. We believe that this is because both rely on PCA to encode blur information. Since the variability of the blur is much higher, it is more difficult to accurately depict new blur kernels at test time. Furthermore, both models rely on standard convolutional layers that are not well suited for recovering blurs caused by strong camera motions. In contrast, the use of the motion-corrected feature extraction sub-module and the Wiener filter allows our model to not only adapt to blur introduced by more
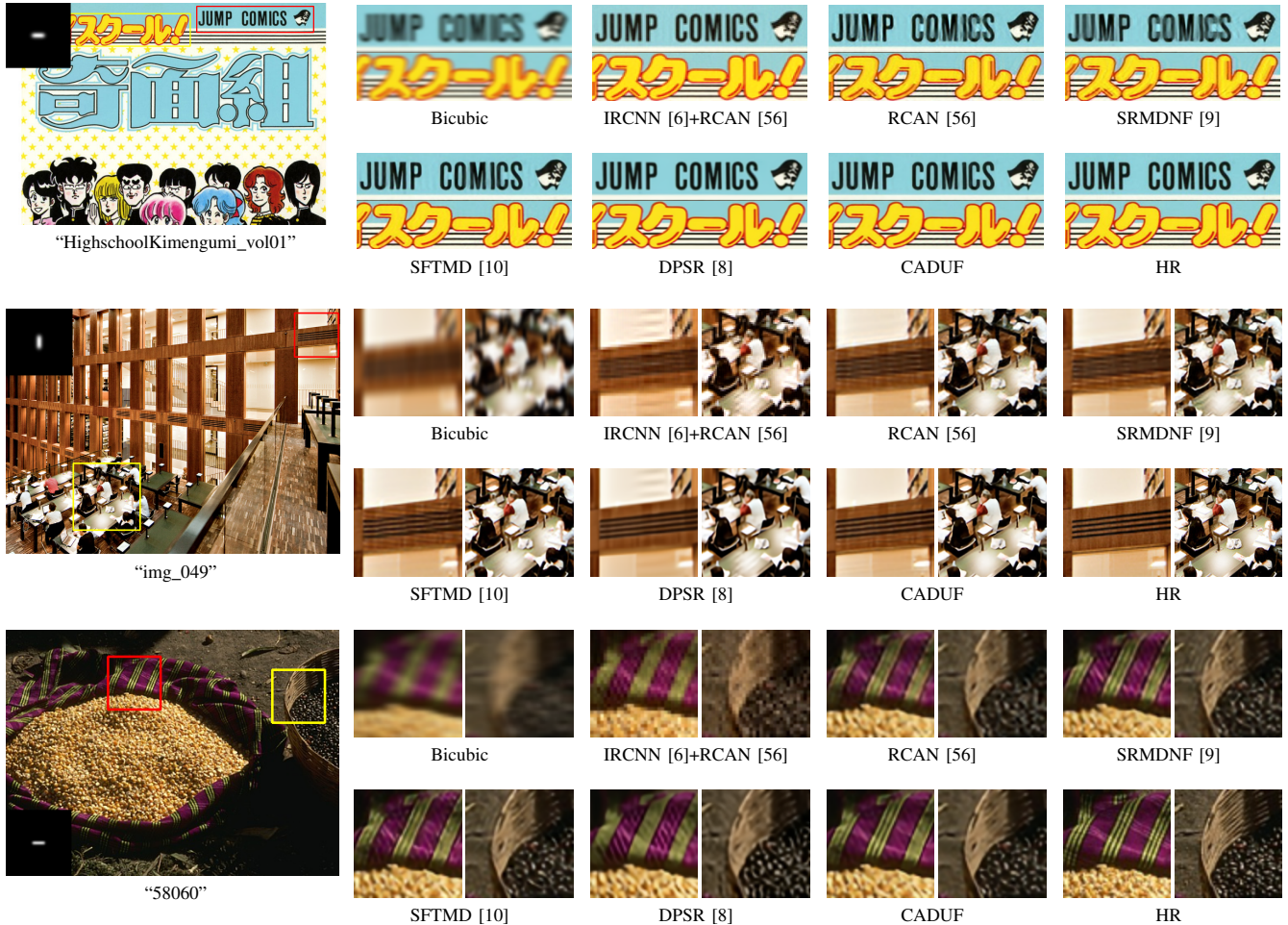
Fig. 6. Comparison of IRCNN [6]+RCAN [56], RCAN [56], SRMDNF [9], SFTMD [10], DPSR [8] and CADUF on GaussianSM dataset for factor 2 on image "HighschoolKimengumi_vol01" from Managa109 [49] dataset, image "img_049" from Urban100 [3] dataset and image "58060" from BSD100 [48] dataset.

complex camera motions but also to better generalize and be more robust to blurs not seen in training. These features are not shared by the other two cascade models IRCNN [6]+RCAN [56] and DeblurGAN+RCAN [56]. Both produce worse results than RCAN [56] and cannot compete with the remaining CNN methods in terms of performance and time consumption. The accumulated errors in the cascade and the strong artifacts that appear due to it are the source of this poor performance.

Compared to MBO methods (ZSSR [34] and DPSR [8]), we can see that CADUF requires orders of magnitude fewer MACs because it does not need to solve an expensive optimization problem for each new test image. Although our model's ability to resolve unseen blurs is lower, it generalizes well to new degradations close to those observed in training, that is, a similar type of noise and blur (motion, defocus, downsampling). We can see that our model significantly outperforms both ZSSR [34] and DPSR [8], with the latter being the closest in terms of performance of the other two. Fig. 7[2] shows a visual comparison of SOTA methods and CADUF for testing images of the GaussianCM dataset. As seen, compared to the other two best-performing models, SFTMD [10] and DPSR [8], CADUF produces significantly better images, closer to the HR image and with far fewer artifacts. Especially remarkable is the case of image "img096" of the Urban100 [3] dataset, where CADUF is able to reconstruct the building windows on the center of the image far better than the other methods. Despite the good performance shown by CADUF, when very strong blurs are present, it would fail to properly recover the HR image. One example of these cases is shown in Fig. 8. The CADUF method is unable to recover the face of the soldier and produces some easily seen artifacts in the background wall and the soldier helmet. However, CADUF still produces better results than the other SOTA methods, introducing significantly fewer artifacts and better recovering the hands of the soldier and more details of the face. When making a comparative visual study of our model against the SOTA, we found that it produces competing results when the blur does not include too much motion and outperforms other models when strong motion is present.

Finally, in Table III, we compare our method against competing methods in the blind setting in the GaussianCM dataset. We estimate the LR blur kernel $k^{*L}$ for algorithms that require

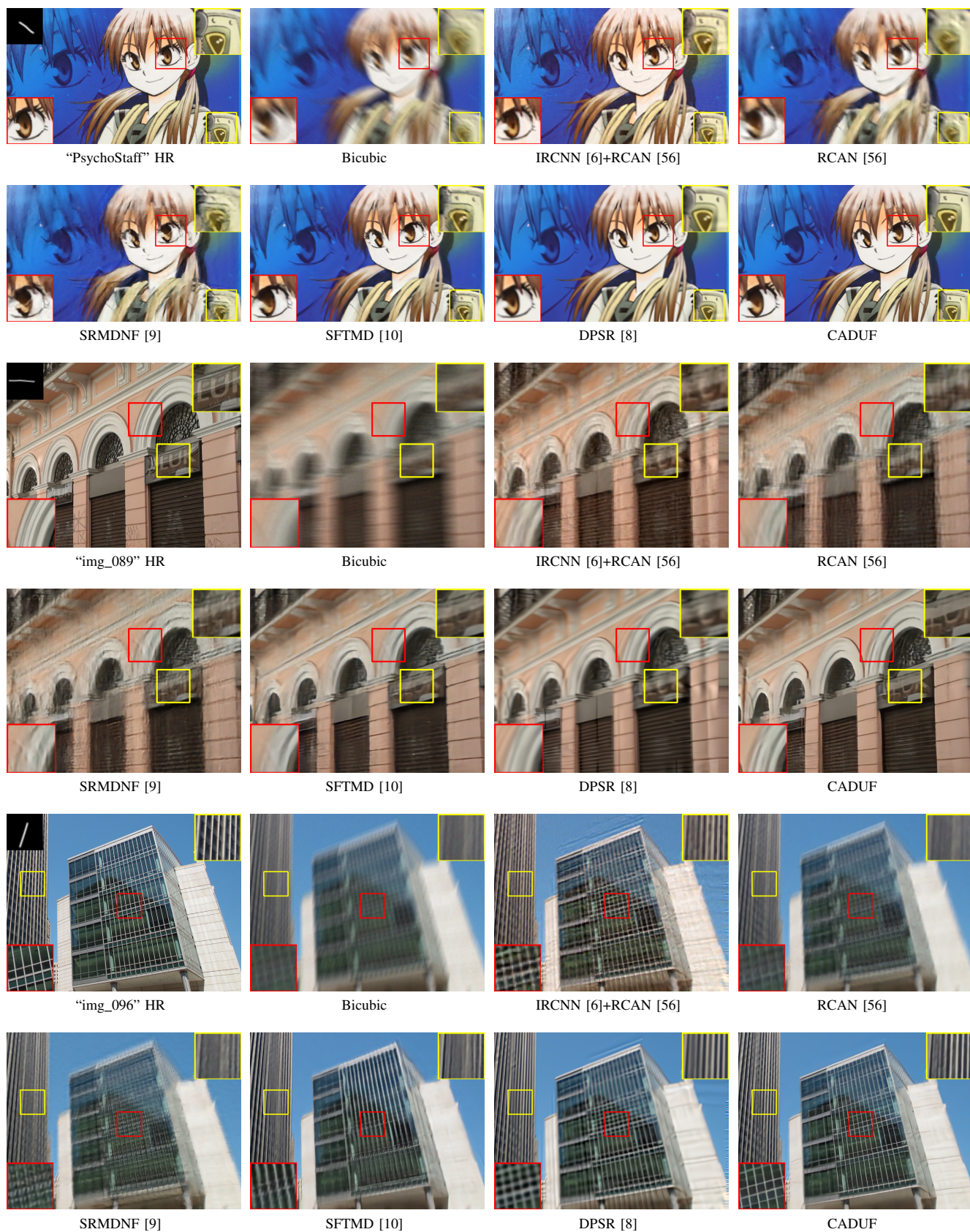[2]More examples can be downloaded at https://github.com/vipgugr/CADUF.

Fig. 7. Visual comparison of IRCNN [6]+RCAN [56], RCAN [56], SRMDNF [9], SFTMD [10], DPSR [8] and CADUF on GaussianCM dataset for factor 2 on image "PsychoStaff" from Managa109 [49] dataset and images "img_089" and "img_096" from Urban100 [3] dataset.
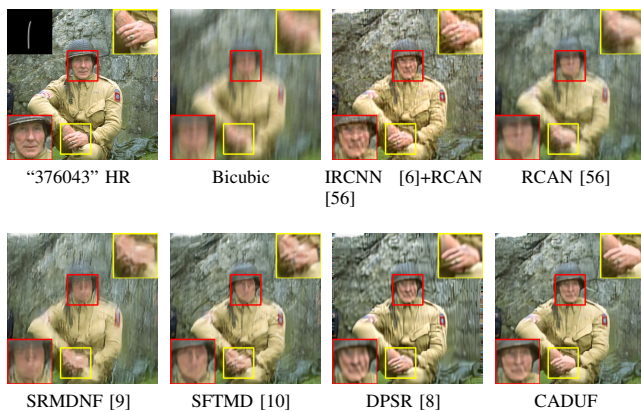
Fig. 8. Failure case of CADUF for factor 2 on image "376043" from the BSD100 [48] dataset with degradation from the GaussianCM dataset.

it (DPSR [8] and ours) using the method in [57]. In our model case, we fine-tuned it for 100 epochs with lr=$10^{-5}$. We report results for factor 2 because the loss of information coupled with errors during the degradation estimation makes it difficult to analyze the differences between the methods for factor 4. As seen, the proposed model has a significantly less drop in performance compared to DPSR [8] when an approximation of $k^{*L}$ is used. In this case, our model benefits of the possibility to fine-tune and adapt to the errors in the blur kernel estimation. Thanks to this fact, it still significantly outperforms all other methods (2.18dB-0.02(SSIM) and 0.61dB-0.06(SSIM) over DPSR [8] and RCAN [56], respectively).

## IV. CONCLUSIONS

In this study, a new cascade CNN model for SISR capable of dealing with multiple degradations, has been proposed. The model exhibits a new approach for searching for a solution with several differentiating keys. First, cascade sub-modules are explicitly specialized in specific tasks: motion-corrected feature extraction, deblur, and upsampling. Second, a new way of regularizing the deblur and upsampling sub-modules is introduced from domain knowledge. As a consequence of the upsampling restriction using affine projection [40], a final sub-module that corrects the accumulated cascade error can be introduced to further increase the model performance. As far we know, the use of domain knowledge to constrain the training of modules in a cascade represents a novelty in SISR. Finally, the model uses an improved method of connecting the sub-modules in the cascade where each sub-module re-uses the features calculated by all previous sub-modules. As a result, our model is fast, robust, and accurate. The experiments were performed using two new SISR datasets, containing degradations with simple and complex camera movement. These experiments show that CADUF significantly outperforms CNN-BMs with architectures of the same depth. Furthermore, CADUF outperforms the remaining compared SOTA methods for non-blind and blind SISR while being significantly faster. Although CADUF is a CNN-BM, it can be generalized to unseen blurs during training, competing in this regard with MBO models.

## REFERENCES

[1] C.-Y. Yang, C. Ma, and M.-H. Yang, "Single-image super-resolution: A benchmark," in *ECCV*, 2014.

[2] C. Dong, C. C. Loy, K. He, and X. Tang, "Learning a deep convolutional network for image super-resolution," in *European Conference on Computer Vision*, pp. 184–199, Springer, 2014.

[3] J. Huang, A. Singh, and N. Ahuja, "Single image super-resolution from transformed self-exemplars," in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 5197–5206, June 2015.

[4] J. Kim, J. K. Lee, and K. M. Lee, "Accurate image super-resolution using very deep convolutional networks," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1646–1654, June 2016.

[5] B. Lim, S. Son, H. Kim, S. Nah, and K. M. Lee, "Enhanced deep residual networks for single image super-resolution," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, vol. 1, p. 3, 2017.

[6] K. Zhang, W. Zuo, S. Gu, and L. Zhang, "Learning deep CNN denoiser prior for image restoration," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2808–2817, July 2017.

[7] X. Wang, K. Yu, C. Dong, and C. Change Loy, "Recovering realistic texture in image super-resolution by deep spatial feature transform," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 606–615, 2018.

[8] K. Zhang, W. Zuo, and L. Zhang, "Deep plug-and-play super-resolution for arbitrary blur kernels," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2019.

[9] K. Zhang, W. Zuo, and L. Zhang, "Learning a single convolutional super-resolution network for multiple degradations," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2018.

[10] J. Gu, H. Lu, W. Zuo, and C. Dong, "Blind super-resolution with iterative kernel correction," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.

[11] A. K. Katsaggelos, R. Molina, and J. Mateos, "Super resolution of images and video," *Synthesis Lectures on Image, Video, and Multimedia Processing*, vol. 1, no. 1, pp. 1–134, 2007.

[12] Y. He, K.-H. Yap, L. Chen, and L.-P. Chau, "A soft map framework for blind super-resolution image reconstruction," *Image and Vision Computing*, vol. 27, no. 4, pp. 364 – 373, 2009.

[13] Y. Tai, S. Liu, M. S. Brown, and S. Lin, "Super resolution using edge prior and single image detail synthesis," in *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pp. 2400–2407, June 2010.

[14] S. P. Belekos, N. P. Galatsanos, and A. K. Katsaggelos, "Maximum a posteriori video super-resolution using a new multichannel image prior," *IEEE Transactions on Image Processing*, vol. 19, no. 6, pp. 1451–1464, 2010.

[15] S. D. Babacan, R. Molina, and A. K. Katsaggelos, "Variational Bayesian super resolution," *IEEE Transactions on Image Processing*, vol. 20, no. 4, pp. 984–999, 2011.

[16] Y. Romano, M. Elad, and P. Milanfar, "The little engine that could: Regularization by denoising (red)," *SIAM J. Imaging Sciences*, vol. 10, pp. 1804–1844, 2017.

[17] D. Geman and C. Yang, "Nonlinear image recovery with half-quadratic regularization," 1995.

[18] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Foundations and Trends in Machine Learning*, vol. 3, pp. 1–122, 2011.

[19] T. Meinhardt, M. Möller, C. Hazirbas, and D. Cremers, "Learning proximal operators: Using denoising networks for regularizing inverse imaging problems," *2017 IEEE International Conference on Computer Vision (ICCV)*, pp. 1799–1808, 2017.

[20] S. A. Bigdeli, M. Zwicker, P. Favaro, and M. Jin, "Deep mean-shift priors for image restoration," in *Advances in Neural Information Processing Systems*, pp. 763–772, 2017.

[21] J. Yang, Z. Wang, Z. Lin, S. Cohen, and T. Huang, "Coupled dictionary training for image super-resolution," *IEEE Transactions on Image Processing*, vol. 21, pp. 3467–3478, Aug 2012.

[22] K. Zhang, X. Gao, D. Tao, and X. Li, "Single image super-resolution with multiscale similarity learning," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 24, pp. 1648–1659, Oct 2013.

[23] D. Trinh, M. Luong, F. Dibos, J. Rocchisani, C. Pham, and T. Q. Nguyen, "Novel example-based method for super-resolution and denoising of medical images," *IEEE Transactions on Image Processing*, vol. 23, pp. 1882–1895, April 2014.

[24] B. C. S. Hui Jung Lee, Dong-Yoon Choi, "Learning-based superresolution algorithm using quantized pattern and bimodal postprocessing for text images," *Journal of Electronic Imaging*, vol. 24, no. 6, pp. 1 – 8 – 8, 2015.

[25] Y. Zhang, Y. Zhang, J. Zhang, D. Xu, Y. Fu, Y. Wang, X. Ji, and Q. Dai, "Collaborative representation cascade for single-image super-resolution," *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, vol. 49, pp. 845–860, May 2019.

[26] C. Dong, C. C. Loy, and X. Tang, "Accelerating the super-resolution convolutional neural network," in *Proceedings of European Conference on Computer Vision*, 2016.

[27] S. Anwar, S. Khan, and N. Barnes, "A deep journey into super-resolution: A survey," *ACM Comput. Surv.*, vol. 53, May 2020.

[28] J. Kim, J. Kwon Lee, and K. Mu Lee, "Accurate image super-resolution using very deep convolutional networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1646–1654, 2016.

[29] C. Ledig, L. Theis, F. Huszár, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang, *et al.*, "Photo-realistic single image super-resolution using a generative adversarial network," *arXiv preprint arXiv:1609.04802*, 2016.

[30] Y. Zhang, Y. Tian, Y. Kong, B. Zhong, and Y. Fu, "Residual dense network for image super-resolution," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2472–2481, June 2018.

[31] Y. Zhang, K. Li, K. Li, L. Wang, B. Zhong, and Y. Fu, "Image super-resolution using very deep residual channel attention networks," in *ECCV*, 2018.

[32] X. Wang, K. Yu, S. Wu, J. Gu, Y. Liu, C. Dong, C. C. Loy, Y. Qiao, and X. Tang, "Esrgan: Enhanced super-resolution generative adversarial networks," in *ECCV Workshops*, 2018.

[33] W. Shi, J. Caballero, F. Huszár, J. Totz, A. P. Aitken, R. Bishop, D. Rueckert, and Z. 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)*, pp. 1874–1883, 2016.

[34] M. I. Assaf Shocher, Nadav Cohen, ""zero-shot" super-resolution using deep internal learning," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.

[35] G. Riegler, S. Schulter, M. Rüther, and H. Bischof, "Conditioned regression models for non-blind single image super-resolution," in *IEEE International Conference on Computer Vision*, pp. 522–530, Dec 2015.

[36] Y. Yuan, S. Liu, J. Zhang, Y. Zhang, C. Dong, and L. Lin, "Unsupervised image super-resolution using cycle-in-cycle generative adversarial networks," *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pp. 814–81409, 2018.

[37] V. Vladimir and I. Rauf, "Learning using privileged information: Similarity control and knowledge transfer," *Journal of Machine Learning Research*, pp. 2023–2049, 2015.

[38] J. Dai, H. Qi, Y. Xiong, Y. Li, G. Zhang, H. Hu, and Y. Wei, "Deformable convolutional networks," *arXiv preprint arXiv:1703.06211*, 2017.

[39] X. Zhu, H. Hu, S. Lin, and J. Dai, "Deformable convnets v2: More deformable, better results," *arXiv preprint arXiv:1811.11168*, 2018.

[40] C. Sonderby, J. Caballero, L. Theis, W. Shi, and F. Huszar, "Amortised MAP inference for image super-resolution," in *International Conference on Learning Representations*, 2017.

[41] A. Albert, *Regression and the Moore-Penrose pseudoinverse*. Mathematics in Science and Engineering, Burlington, MA: Elsevier, 1972.

[42] E. Pérez-Pellitero, M. S. M. Sajjadi, M. Hirsch, and B. Schölkopf, "Photorealistic video super resolution," in *Workshop and Challenge on Perceptual Image Restoration and Manipulation (PIRM) at the 15th European Conference on Computer Vision (ECCV)*, 2018.

[43] F. Yu and V. Koltun, "Multi-scale context aggregation by dilated convolutions," in *International Conference on Learning Representations (ICLR)*, May 2016.

[44] X. Jia, B. De Brabandere, T. Tuytelaars, and L. V. Gool, "Dynamic filter networks," in *Advances in Neural Information Processing Systems 29* (D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, eds.), pp. 667–675, Curran Associates, Inc., 2016.

[45] Y. Jo, S. W. Oh, J. Kang, and S. J. Kim, "Deep video super-resolution network using dynamic upsampling filters without explicit motion compensation," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3224–3232, 2018.

[46] K. Zhang, W. Zuo, Y. Chen, D. Meng, and L. Zhang, "Beyond a Gaussian denoiser: Residual learning of deep CNN for image denoising," *IEEE Transactions on Image Processing*, vol. 26, no. 7, pp. 3142–3155, 2017.

[47] D. Li and Z. Wang, "Video super-resolution via motion compensation and deep residual learning," *IEEE Transactions on Computational Imaging*, 2017.

[48] D. Martin, C. Fowlkes, D. Tal, and J. Malik, "A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics," in *Proceedings Eighth IEEE International Conference on Computer Vision. ICCV 2001*, vol. 2, pp. 416–423 vol.2, July 2001.

[49] Y. Matsui, K. Ito, Y. Aramaki, A. Fujimoto, T. Ogawa, T. Yamasaki, and K. Aizawa, "Sketch-based manga retrieval using manga109 dataset," *Multimedia Tools and Applications*, vol. 76, pp. 21811–21838, Oct 2017.

[50] E. Agustsson and R. Timofte, "Ntire 2017 challenge on single image super-resolution: Dataset and study," in *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pp. 1122–1131, July 2017.

[51] R. Timofte, E. Agustsson, L. V. Gool, M. Yang, L. Zhang, B. Lim, S. Son, H. Kim, S. Nah, K. M. Lee, X. Wang, Y. Tian, K. Yu, Y. Zhang, S. Wu, C. Dong, L. Lin, Y. Qiao, C. C. Loy, W. Bae, J. Yoo, Y. Han, J. C. Ye, J. Choi, M. Kim, Y. Fan, J. Yu, W. Han, D. Liu, H. Yu, Z. Wang, H. Shi, X. Wang, T. S. Huang, Y. Chen, K. Zhang, W. Zuo, Z. Tang, L. Luo, S. Li, M. Fu, L. Cao, W. Heng, G. Bui, T. Le, Y. Duan, D. Tao, R. Wang, X. Lin, J. Pang, J. Xu, Y. Zhao, X. Xu, J. Pan, D. Sun, Y. Zhang, X. Song, Y. Dai, X. Qin, X. Huynh, T. Guo, H. S. Mousavi, T. H. Vu, V. Monga, C. Cruz, K. Egiazarian, V. Katkovnik, R. Mehta, A. K. Jain, A. Agarwalla, C. V. S. Praveen, R. Zhou, H. Wen, C. Zhu, Z. Xia, Z. Wang, and Q. Guo, "Ntire 2017 challenge on single image super-resolution: Methods and results," in *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pp. 1110–1121, July 2017.

[52] G. Boracchi and A. Foi, "Modeling the performance of image restoration from motion blur," *IEEE Transactions on Image Processing*, vol. 21, no. 8, pp. 3502–3517, 2012.

[53] W.-S. Lai, J.-B. Huang, N. Ahuja, and M.-H. Yang, "Fast and accurate image super-resolution with deep laplacian pyramid networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018.

[54] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, "Pytorch: An imperative style, high-performance deep learning library," in *Advances in Neural Information Processing Systems 32* (H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, eds.), pp. 8024–8035, Curran Associates, Inc., 2019.

[55] D. Ulyanov, A. Vedaldi, and V. S. Lempitsky, "Deep image prior," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

[56] Y. Zhang, K. Li, K. Li, L. Wang, B. Zhong, and Y. Fu, "Image super-resolution using very deep residual channel attention networks," *CoRR*, vol. abs/1807.02758, 2018.

[57] X. Zhou, M. Vega, F. Zhou, R. Molina, and A. K. Katsaggelos, "Fast Bayesian blind deconvolution with Huber super Gaussian priors," *Digital Signal Processing*, vol. 60, pp. 122–133, 2017.