# 3DSRnet: Video Super-resolution using 3D Convolutional Neural Networks

| Soo Ye Kim | Jeongyeon Lim | Taeyoung Na | Munchurl Kim |
|:---:|:---:|:---:|:---:|
| KAIST | SK Telecom | SK Telecom | KAIST |
| sooyekim@kaist.ac.kr | jeongyeon@sk.com | taeyoung.na@sk.com | mkimee@kaist.ac.kr |

## Abstract

*In video super-resolution, the spatio-temporal coherence between, and among the frames must be exploited appropriately for accurate prediction of the high resolution frames. Although 2D convolutional neural networks (CNNs) are powerful in modelling images, 3D-CNNs are more suitable for spatio-temporal feature extraction as they can preserve temporal information. To this end, we propose an effective 3D-CNN for video super-resolution, called the 3DSRnet that does not require motion alignment as preprocessing. Our 3DSRnet maintains the temporal depth of spatio-temporal feature maps to maximally capture the temporally nonlinear characteristics between low and high resolution frames, and adopts residual learning in conjunction with the sub-pixel outputs. It outperforms the most state-of-the-art method with average 0.45 and 0.36 dB higher in PSNR for scales 3 and 4, respectively, in the Vidset4 benchmark. Our 3DSRnet first deals with the performance drop due to scene change, which is important in practice but has not been previously considered.*

## 1. Introduction

Vision is one of the most primitive yet sophisticated sensory systems that is continuously stimulated not only by natural scenes, but also by electric displays. With the unceasingly evolving display hardware which has now commercially reached the resolution of 8K Ultra High Definition (UHD), and people's rising expectations on these types of visuals, the demand for better quality videos is at its highest. However, the sole advancement of display technologies is not sufficient to offer high quality visual content to users - the contents themselves have to be of higher resolution. Although they can be obtained through the usage of high-end filming equipment, it is costly and problematic due to large storage and transmission bandwidth required.

Super-resolution (SR) is an imaging technique that transforms low resolution (LR) images to higher resolution ones. When an LR image is given as input, an SR algorithm exploits its internal information to generate an output image,
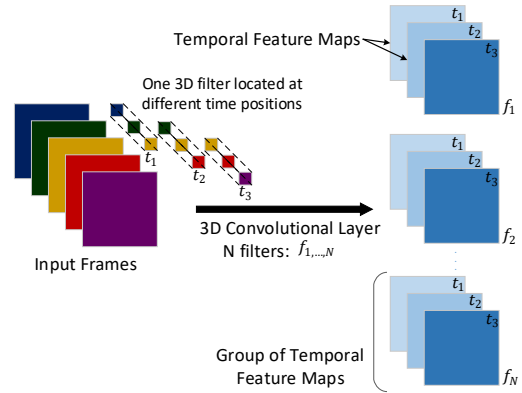


Figure 1. Illustration of a 3D convolution layer with a five-frame input and filters of depth 3. When *N* 3D filters are applied to the input, each of the filters generate a group of temporal feature maps, resulting in *N* groups of temporal feature maps

hopefully similar to its high resolution (HR) counterpart. This is regarded as an ill-posed problem since multiple HR images correspond to a single LR image. Non-existent, but reasonable, information should be created within the image when going from LR to HR, and finding a high quality image among the possible solutions is the key to the SR problem.

Despite that SR is a popular problem in image processing and computer vision, most studies have focused on single image SR than multi-frame SR, also referred to as video SR. However, many SR applications are in videos where the reconstruction of HR frames may benefit from additional information contained in the previous and future LR frames. While video frames exhibit high temporal coherence, the camera or object motion can also provide a different angle or scale of the parts in the current frame in the consecutive surrounding frames, which can be effectively utilized as crucial clues in constructing high quality HR frames.

A video SR algorithm should fully exploit the temporal relations between the consecutive frames to aggregate them with the spatial information. To this end, we propose an effective 3D convolutional neural network (CNN) for video SR, called 3DSRnet that does not require motion estima-

tion nor compensation to interpret the spatio-temporal information in consecutive frames. Instead, it finds an end-to-end nonlinear spatio-temporal mapping in itself through residual learning and lowers complexity by using the multi-channel output structure introduced in [1]. Our 3DSRnet outperforms the previous video SR methods [2, 3] by at least average 0.36 dB in PSNR for the Vidset4 benchmark test dataset. To the best of our knowledge, it is also the first video SR method that can effectively deal with scene change in the input frames.

## 2. Related Work

### 2.1. Single Image Super-resolution

Single image SR attempts to develop an HR image from a single LR image. Past attempts to tackle this problem include internal and external example-based methods [4, 5, 6, 7, 8]. The former includes a method devised by Glasner et al. [6], which identifies internal redundancies in an image to obtain essential information in upscaling of the patches. The external example-based methods try to find a dictionary mapping [4, 5, 8]. Another type of approach is through sparse representation, applied successfully by Yang et al [7].

With the recent rise of deep learning and the excellent performance of CNNs in image classification [9], the first structure that adopts a CNN structure for SR was proposed by Dong et al. [11, 10], which suggests a simple 3-layer structure. Their model, called SRCNN, demonstrated great potential of using CNNs for SR applications. Since then, CNN-based structures have been boasting superior performance. One of the CNN-based SR methods that was highly successful is called the very deep super-resolution method (VDSR) proposed by Kim et al [12]. The VDSR has as many as twenty convolution layers and first adopts residual learning to train a deep SR network. However, both SRCNN and VDSR start with enlarged LR images using a bicubic filter, as input to the first convolution layer. Consequently, the convolution operations are taken place on the enlarged input, which leads to high computation complexity. An inspirational work by Shi et al. [1] suggested a sub-pixel CNN that finds a direct transform from the LR image by using the fact that convolution layers can produce multiple channels at the output. With this multi-channel output structure, the HR image can be obtained through a simple reordering of the output pixels. Our 3DSRnet employs this multi-channel output structure [1] with residual learning in [12].

### 2.2. Video Super-resolution

Video SR, or multi-frame SR, assumes that the input is a series of consecutive frames at each time instance of video sequences. Undoubtedly, single image SR algorithms may be applied on the individual frames for videos, and this may even be more efficient in some cases if they achieve real-time performance as in [1]. However, more spatial information is available in the case of videos, as not only the current LR frame but also its surrounding consecutive LR frames may be utilized. This means that to fully profit from what is given, the temporal relation of the spatial information has to be carefully taken into account in reconstructing the corresponding HR frame.

Compared to image SR, relatively less studies have been conducted on video SR. Focusing on neural network based methods, Kappeler et al. [3] extended SRCNN to 2D-CNN architectures that combine information from neighboring frames. Caballero et al. [2] proposed three video SR architectures where the early and slow fusion architectures have a similar way of dealing with the multi-frame input as in [3, 13]. The structures in [2] all adopt the same multi-channel output structure in [1]. The third model is a 3D-CNN architecture that first incorporated 3D convolution filters into video SR to capture temporal information of multiple frames. This model was a conceptual suggestion without the specific configuration information presented, where the temporal depth of feature maps shrinks to one in early convolution layers on which 2D convolutions are then performed. No performance comparison for the 3D-CNN architecture [2] was provided against other previous methods due to its relatively lower SR performance compared to the early and slow fusion architectures.

In comparison, our 3DSRnet maintains the temporal depth of spatio-temporal feature maps towards deeper layers to maximally capture the temporally nonlinear characteristics between LR and HR frames, and we provide intensive experiments and analysis on 3DSRnet in the later sections of this paper. All three video SR architectures in [2] and the video SR method in [3] need motion alignment among the multiple input frames whereas our 3DSRnet directly takes the input frames without any motion alignment. None of the previous CNN-based video SR methods has considered scene change issues while our 3DSRnet first incorporates a scene change detection network to locate a scene change boundary in multiple input frames and replace the different scene frames with the temporally closest frame of the same scene as the current frame scene.

### 2.3. 3D Convolutional Neural Networks

2D-CNN, a common neural network used for images, is a powerful structure in modelling images with their spatial feature extraction capability. However when another axis, time, is introduced as in videos, we argue that 3D-CNN is a more suitable option for spatio-temporal feature extraction. This is in line with Tran et al. [14] in which they argued that 3D-CNN is an effective video descriptor and with [15] that demonstrated the spatial and temporal feature extraction capability of 3D-CNNs. 3D-CNNs have been success-
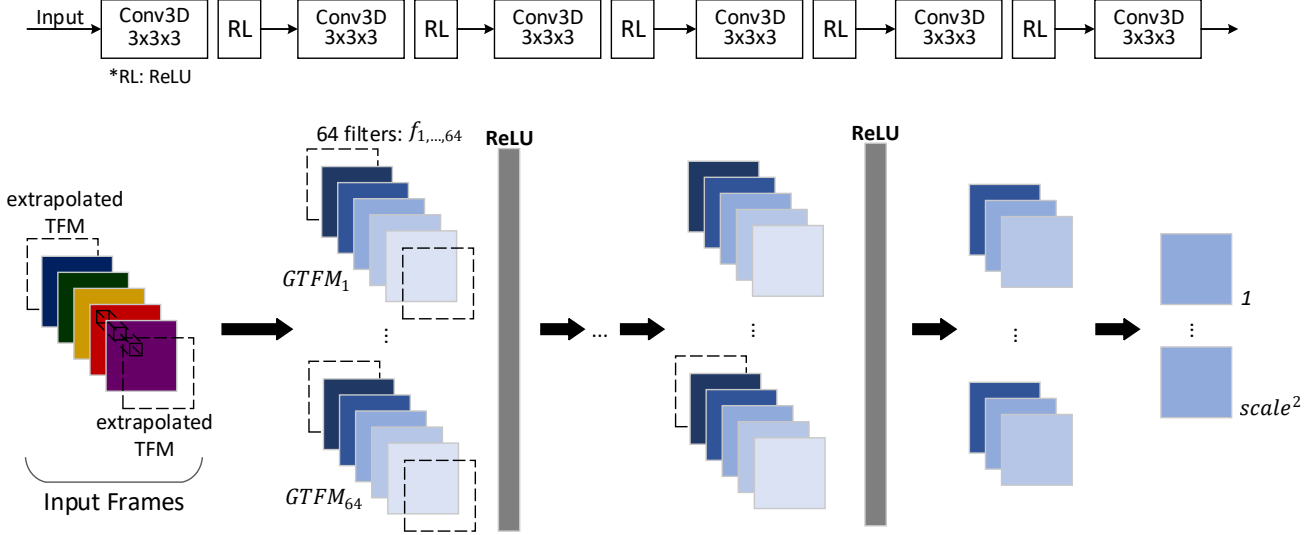
Figure 2. Architecture of the video SR subnet in 3DSRnet. Each convolution layer has 64 filters of size 3×3×3 and is followed by ReLU activation except the last layer. Each filter produces a group of temporal feature maps (GTFM) and a temporal feature map is extrapolated on both ends of each GTFM to preserve the temporal depth which would otherwise decrease. No extrapolation is performed from layer *L-1* so that the temporal information is merged to produce the final 2D output. The temporal information is preserved until the last layer

fully implemented in high level vision tasks for videos such as action/object recognition and scene/event classification [14, 15, 16]. We believe they are also effectively applicable to a low level vision task for videos such as video SR. In this paper, we adopt the 3D-CNN and design an elaborate video SR network, 3DSRnet, which makes the 3D convolution effective on video SR where motion alignment is not necessitated thanks to its spatio-temporal feature representation ability.

# 3. Proposed Method

We propose a 3D-CNN architecture for video SR named as the 3DSRnet with an additional scene change module that deals with scene change occurring inputs. The 3DSRnet consists of two subnets:

(i) Video SR subnet, and

(ii) Scene change detection and frame replacement (SF) subnet.

The video SR subnet takes a series of consecutive LR input frames in a sliding time window, and produces an HR output frame corresponding to the middle frame in the sliding time window. The SF subnet of the 3DSRnet is responsible for the detection of scene change in the sliding time window, and replaces the frames of a different scene with the temporally closest frame that belongs to the same scene as the middle frame.

## 3.1. Video Super-resolution Subnet

### 3.1.1 3D Convolution Layers.

The video SR subnet is composed of 3D convolution layers where 3D filters of size $height \times width \times depth$ are applied on the input composed of multiple consecutive frames or feature maps. Unlike 2D filters of size $height \times width$ that are applied on the full depth of the input and slid horizontally and vertically, 3D filters have a third size parameter, $depth$, so that they are swept horizontally, vertically and depth-wise. The first 3D convolution layer takes a series of five consecutive frames in a sliding input window where each 3D filter generates a temporal feature map (TFM) at the corresponding frame position, each filter yielding a group of temporal feature maps (GTFMs) from all the input frame positions. This is illustrated in Fig. 1 where $N$ 3D filters, $f_{1,...,N}$ of depth 3 are applied on the input composed of five frames to produce $N$ GTFMs. The temporal depth of each GTFM is 3. Formally, the $n$-th GTFM before activation in the first 3D convolution layer is given by

$$GTFM^n_{xyt} = \sum_{h=1}^{H} \sum_{w=1}^{W} \sum_{d=1}^{D} w^n_{hwd} v_{(x+h)(y+w)(t+d)} + b^n$$

$$(1)$$

where $w^n_{HWD}$ is the 3D filter *n* of size $H \times W \times D$ and *v* is the input.

From the second to the last 3D convolution layer, the input window corresponds to the whole set of multiple GTFMs where each GTFM is generated from one 3D filter of the previous convolution layer. The temporal information
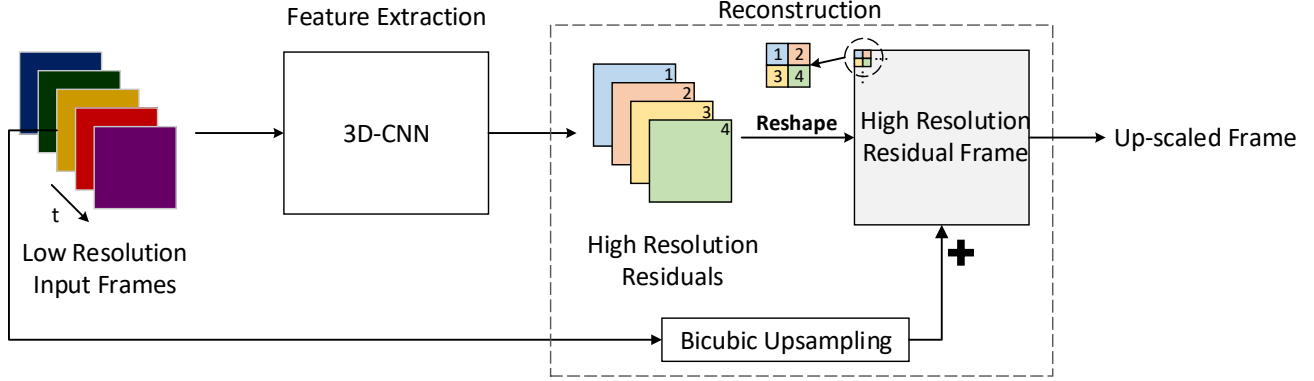
Figure 3. The input and output structure of 3DSRnet. This is an example with five input frames and the scale factor of 2. The input frames go through a 3D-CNN for spatio-temporal feature extraction and this 3D-CNN predicts the high resolution residuals of the middle frame. There are four output channels because the scale factor is 2 (*no. of output channels=$scale^2$*). They are reshaped to make the high resolution residual frame. The final up-scaled frame is obtained by adding the bicubic up-scaled middle frame to the residual frame

contained within the input window is preserved through the 3D convolution layer in each GTFM as separate TFMs unlike 2D convolution layers where the input would be collapsed into one single feature map per filter. From the second to the last 3D convolution layer, the input *v* is composed of multiple GTFMs (*m* GTFMs), and the *n*-th GTFM before activation is given by

$$GTFM_{xyt}^n = \sum_m \sum_{h=1}^H \sum_{w=1}^W \sum_{d=1}^D w_{hwd}^{mn} v_{(x+h)(y+w)(t+d)}^m + b^n.$$
(2)

We use the ReLU [17] function as the activation function after every convolution layer except the last, for the nonlinearity of the network. In 3D-CNNs, temporal nonlinearities among the GTFMs as well as spatial nonlinearities are introduced thanks to the 3D convolution structure.

### 3.1.2 Extrapolation.

The temporal depth of GTFMs become shallower as the network gets deeper, as the 3D filters integrate the temporal information. For example, with an input of five frames and a 3D filter of depth 3, the output GTFMs would have depth 1 after only two 3D convolution layers, being no different from a 2D convolution layer from layers thereafter. Since the usage of 3D-CNNs is to tamper with the temporal information, thereby introducing temporal nonlinearities, extrapolating (or padding) the input GTFMs at their front and back ends allows to preserve the temporal depth throughout the network. However, for the last layers, no extrapolation is performed and the temporal information is aggregated, to produce the final 2D HR frame as intended. For an input of five frames and a 3D filter of depth 3, no extrapolation is carried out from layer *L-1* where *L* is the number of convolution layers. This naturally aggregates the temporal information when going deeper in the network. Please refer to Fig. 2 for a detailed illustration of the 3D convolution layers with extrapolation.

### 3.1.3 Multi-channel Output.

The multi-channel output structure first introduced in [1] allows for a direct mapping from the LR to HR frames by producing an output with multiple channels that can simply be reordered and reshaped to produce the final HR output. This method alleviates the amount of computation which can be otherwise expensive for 3D-CNNs. Furthermore, it can enhance SR performance because the receptive field of the LR input pixels without bicubic up-scaling is larger than that of an up-scaled LR input pixels, provided that the filter size and network depth are the same. Large receptive fields are essential in SR to yield high performance [12, 18, 19, 20].

### 3.1.4 Residual Learning.

HR frames often consist of low and high frequency components. However, the low frequency components are mostly present in the LR frames, meaning that the essential goal of an SR algorithm is in predicting the missing high frequency components. Therefore, the network can save the trouble of predicting what is already there, by directly predicting the difference between the HR frame and the corresponding bicubic-upscaled LR frame - the *residual* frame -. Our 3DSRnet employs this technique and predicts the residual frame, producing a *multi-channel residual* output. Residual learning was first proposed in [21] and applied to SR in [12]. It also eases training [12] by solving the vanishing and exploding gradient problem which can be critical in training neural networks [22]. Fig. 3 shows the input and output structures of the 3DSRnet.
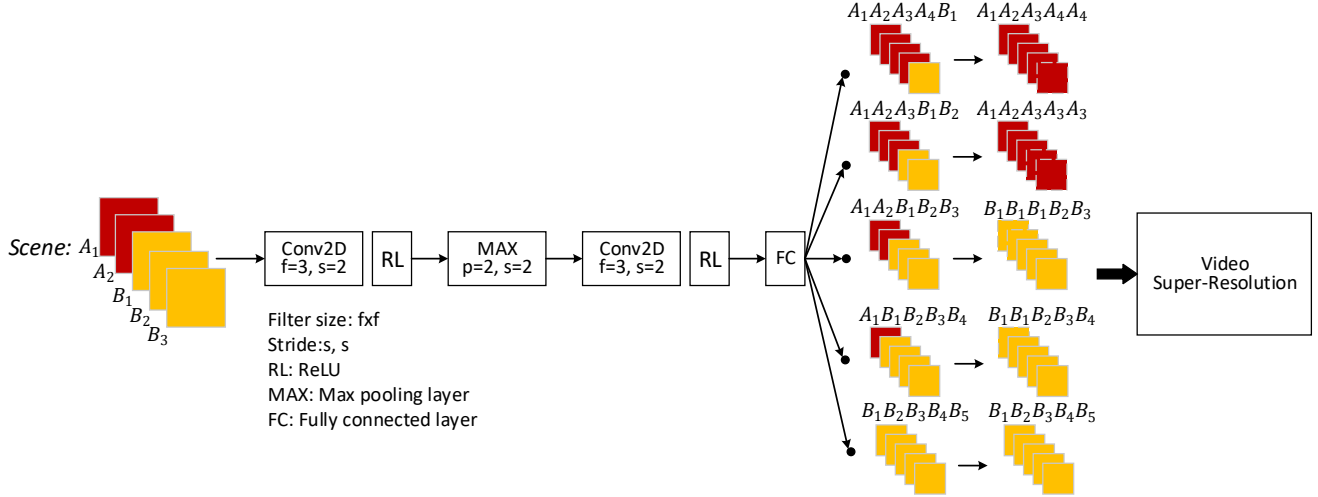
4

Figure 4. Illustration of the 3-layer SF subnet for a sliding time window with five consecutive input frames. The squares in red and yellow colors denote the frames of scene A and scene B, respectively. The input is classified into whether a scene change occurs after frame index 1, 2, 3, 4 or not. After classification, the frames of a different scene are replaced with the temporally closest frame that belongs to the same scene as the middle frame. Note that the middle frame is considered as the reference, and the frames of a different scene than the middle frame are swapped with that of the same scene

## 3.2. Scene Change Detection and Frame Replacement Subnet

The scene change detection and frame replacement (SF) subnet is another component of the 3DSRnet. When multiple frames are used as input to video SR networks, there is a possibility of scene change within them. In this case, the performance of a video SR algorithm drops due to the frames of different scenes getting involved into convolution, resulting in the reconstructed HR frames of poor quality. The previous video SR methods avoided this problem by explicitly collecting data without scene changes, which is impractical in real world applications. Our 3DSRnet handles the scene change problem for video SR by introducing the SF subnet that classifies the exact location of the scene boundary and modifies some of the frames in the sliding input window by replacing the different scene frames with the temporally closest frame of the same scene as the current frame scene. Although the duplicated (replaced) frames do not contain any new information, this method significantly helps the 3DSRnet alleviate performance degradation from the contaminated input of different scene frames.

### 3.2.1 Identification of Scene Change Location.

If we assume that a scene change may occur within the five input frames, there are four possible scene change locations (labels). In addition, the fifth label is designated for no scene change. Then it is a simple five-class classification problem. Fig. 4 illustrates the detailed mechanism of the SF subnet for a sliding time window of five consecutive input frames. The SF subnet should be lightweight as it can

be optionally used alongside video SR, but accurate to correctly modify the input. Therefore, we use a shallow 2D-CNN structure. It is trained separately from the video SR subnet.

## 4. Experiments

### 4.1. Experiment Conditions

#### 4.1.1 Data.

A training or testing data sample of 3DSRnet is composed of five bicubic-down-scaled LR frames and a single HR middle frame. We collected two sets of $3840\times2160$ UHD videos of 30 fps that were encoded with at least 100 Mb/s using an H.264/AVC encoder. The first video set (Type 1) shows spatially complex scenes, meaning that they contain sophisticated objects such as the bird view of a city, and the second video set (Type 2) is temporally complex, meaning that there is a lot of motion. We collected three Type 1 videos of total 8,504 frames and one Type 2 video of 8,655 frames. They were converted into 420 YUV format and only the Y channel was used as the training and test data.

Table 1. Training data sets for the video SR subnet

| Dataset | Type 1 | | Type 2 | | Total no. of subim. |
| --- | --- | --- | --- | --- | --- |
| | stride | subim./fr. | stride | subim./fr. | |
| *smallSet* | 30 | 5 | 5 | 10 | 9,725 |
| *largeSet* | 15 | 10 | 5 | 12 | 14418 |

*subimages and frame are shortened as subim and fr, respectively.

5

Table 2. Experiment on architectures trained on the *smallSet*

| Layers | 2D-CNN | | 3DSRnet v1 | | 3DSRnet v2 | | 3DSRnet v3 | | 3DSRnet | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Number of filter channels (input, output) | | | | | | | | | |
| 1 | 2D | 5, 32 | 3D | 1, 32 | 3D | 1, 32 | 3D | 1, 32 | 3D | 1, 32 |
| 2 | 2D | 32, 64 | 3D | 32, 32 | 3D | 32, 32 | 3D | 32, 32 | 3D | 32, 32 |
| 3 | 2D | 64, 64 | 3D | 32, 16 | 3D | 32, 32 | 3D | 32, 32 | 3D | 32, 32 |
| 4 | 2D | 64, 64 | 2D | 80, 64 | 3D | 32, 16 | 3D | 32, 32 | 3D | 32, 32 |
| 5 | 2D | 64, 35 | 2D | 64, 32 | 2D | 80, 64 | 3D | 32, 16 | 3D | 32, 32 |
| 6 | 2D | 35, 4 | 2D | 32, 4 | 2D | 64, 4 | 2D | 80, 4 | 2D | 32, 4 |
| 2D filter size | 3×3 | | 3×3 | | 3×3 | | 3×3 | | 3×3 | |
| 3D filter size | - | | 3×3×3 | | 3×3×3 | | 3×3×3 | | 3×3×3 | |
| Concat. layer | - | | 3 | | 4 | | 5 | | - | |
| Total parameters | 115,020 | | 108,000 | | 118,368 | | 100,512 | | 114,912 | |
| PSNR (dB) | 32.49 | | 32.81 | | 32.89 | | 32.85 | | **32.92** | |

Table 3. Vidset4 Benchmark-Video SR Methods

| Vidset4 | Bayesian [26] | Deep-DE [27] | VSRnet [3] | Liu *et al.* [28] | VESPCN [2] | 3DSRnet |
|---|---|---|---|---|---|---|
| ×4 | 24.66 | 24.68 | 24.84 | 25.24 | 25.35 | **25.71** |

When reconstructing color frames, U and V channels were simply up-scaled using a bicubic filter.

For training the video SR subnet in 3DSRnet, we prepared two datasets, *smallSet* and *largeSet*, where a predefined number of non-overlapping subimages were randomly selected from each frame with a frame stride from Type 1 and Type 2 sets. For fair comparison with other video SR methods, the video SR subnet was trained with a training dataset without scene change. Table 1 summarizes the training sets for the video SR subnet of the 3DSRnet. The size of LR subimages for the scale factors 2, 3 and 4 were 80, 60, 40 for the *smallSet* and 80, 60, 45 for the *largeSet*, respectively. The training took around three days with the *smallSet* and eight days with the *largeSet* using an Nvidia TITAN X GPU for a scale factor of 2. For the comparison among the video SR subnet of 3DSRnet and its variants, the test set contains the data samples of scenes that are not included in the training set. To compare the video SR subnet of 3DSRnet with the state-of-the-art SR methods, we used the Vidset4 dataset which is a commonly used test set for videos.

For training the SF subnet in 3DSRnet, a separate dataset was created to contain scene changes. The LR frames of different scenes from the smallSet were reduced by a factor of 40 to be of size 48×27, and randomly concatenated to make scene change occurring inputs. 2,000 data samples each consisting of the frames and its label were randomly selected for each of the five classes to make the final training data of 10,000 data samples.

#### 4.1.2 Training.

The purpose of the video SR subnet was to minimize the mean squared loss between the predicted frame $F(X_i; \theta)$ and the ground truth frame $Y_i$, given by

$$Loss(\theta) = \frac{1}{2n} \sum_{i=1}^{n} \|F(X_i; \theta) - Y_i\|^2 \qquad (3)$$

where $X_i$ is the input frames, $\theta$ is the set of model parameters and $n$ is the number of data samples. Then the gradient is calculated as the difference between $F(X_i; \theta)$ and $Y_i$. All weights were initialized by the Xavier initialization [23] using both the number of input and output neurons of the layer. The parameters were updated using Adam [24].

All 3D filter sizes were empirically set to 3×3×3, 2D filters to size 3×3 and the number of filters are 64 if not otherwise mentioned. The network is composed of six convolution layers, considering the tradeoff between performance and complexity. The learning rate was set to $5\times10^{-4}$ for the *smallSet* and $10^{-4}$ for the *largeSet*. The learning rates of biases are 10 times smaller. For all network models, the weight decay was set to $5\times10^{-4}$ for filters and zero for biases. The mini-batch size is 32 for the *smallSet* and 64 for the *largeSet*. All models were implemented using the MatConvNet [25] package and 3D convolution layers were added using a Matlab mex implementation available in GitHub[1].

---

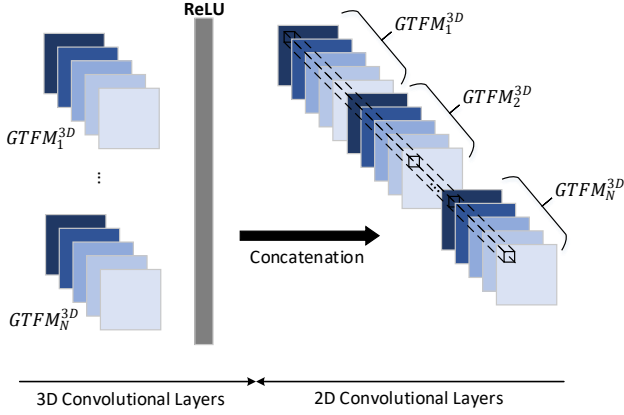[1] https://github.com/pengsun/MexConv3D

6

Figure 5. Concatenation layer for 3D-2D integrated architectures. The GTFMs generated from the last 3D convolution layer are concatenated as input for the following first 2D convolution layer

## 4.2. Architecture

### 4.2.1 Effect of 3D-CNNs over 2D-CNNs.

The video SR subnet in 3DSRnet takes a 3D input (multiple LR input frames) in a sliding time window at a time instance and produces one single 2D HR output frame. So its architecture must be devised to go from 3D to 2D. As illustrated in in Fig. 2, the temporal depth of the GTFMs is kept constant until the (*L-2*)-th convolution layer, and from the (*L-1*)-th convolution layer, no more temporal extrapolation is done to gradually reduce the temporal depth of GTFMs to 1 for our 3DSRnet. As variants of the 3DSRnet, we also experimented with the combination of 3D and 2D convolution layers by simply concatenating the GTFMs created from the last 3D convolution layer and performing 2D filtering thenceforth, with the number of filters adjusted so that all architectures have a similar number of parameters. The concatenation layer is illustrated in Fig. 5. Table 2 summarizes the specifications and results of our 3DSRnet and its variants, 3DSRnet v1, 3DSRnet v2 and 3DSRnet v3 with the comparison to a 2D-CNN structure, also with a five-frame input, made available to demonstrate the superior feature extraction capability of 3D-CNNs.

### 4.2.2 Feature map visualizations.

Fig. 6 shows the feature maps produced from the first convolution layer of the 2D-CNN and 3DSRnet, experimented in Table 2. Feature maps of 3DSRnet appears much sharper, due to the shorter time window length of three. 2D-CNN convolves all five frames at the first convolution layer, producing more blurry feature maps. Furthermore, a 3D filter in 3DSRnet produces a GTFM containing five TFMs for each time instant, preserving the temporal information.
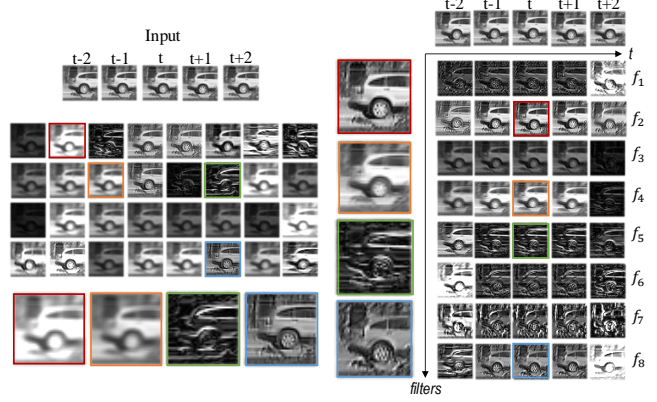


(a) Feature Maps of 2D-CNN     (b) Feature Maps of 3DSRnet

Figure 6. Feature map visualization of 2D-CNN and 3DSRnet

### 4.2.3 Extrapolation.

The video SR subnet of our 3DSRnet extrapolates the GTFMs with a TFM on their both ends to preserve the temporal depth towards the network's deeper layers. There are different ways of extrapolation such as simply padding them with TFMs filled with zeros or with the outmost TFMs duplicated. However, empirical results showed that there was an insignificant difference in performance with 32.88 dB for duplicate extrapolation and 32.92 dB for zero-filled extrapolation. For simplicity, we choose to use the zero-filled extrapolation.

## 4.3. Benchmark.

We test the 3DSRnet for quantitative evaluation in comparison with the state-of-the-art image and video SR methods for the Vidset4 dataset - a popular benchmarking test set that contains four video sequences, namely *Calendar*, *City*, *Foliage* and *Walk*. The PSNR comparison against other video SR methods on scale 4 is given in Table 3, and the PSNR and SSIM comparison against image and video SR methods on scale factors 2, 3 and 4 are given in Table 5. The results of video SR methods [2, 3] are the reported performance on the same test set. The results of [26], [27] and [28] are from those reported in [28]. The image SR methods [11, 12] were tested on the set using their respective codes provided by the authors. As shown in Table 3 and 5, our 3DSRnet outperforms all the state-of-the-art image and video SR methods. Note that in Table 5, the 3DSRnet

Table 4. PSNR (dB) results on each sequence of Vidset4

| Vidset 4 | *Calendar* | *City* | *Foliage* | *Walk* |
|---|---|---|---|---|
| ×2 | 27.04 | 34.13 | 31.58 | 36.25 |
| ×3 | 24.19 | 28.25 | 27.42 | 30.95 |
| ×4 | 22.41 | 26.81 | 25.23 | 28.38 |

*trained with the *largeSet*.

Table 5. Vidset4 Benchmark-Image and Video SR Methods

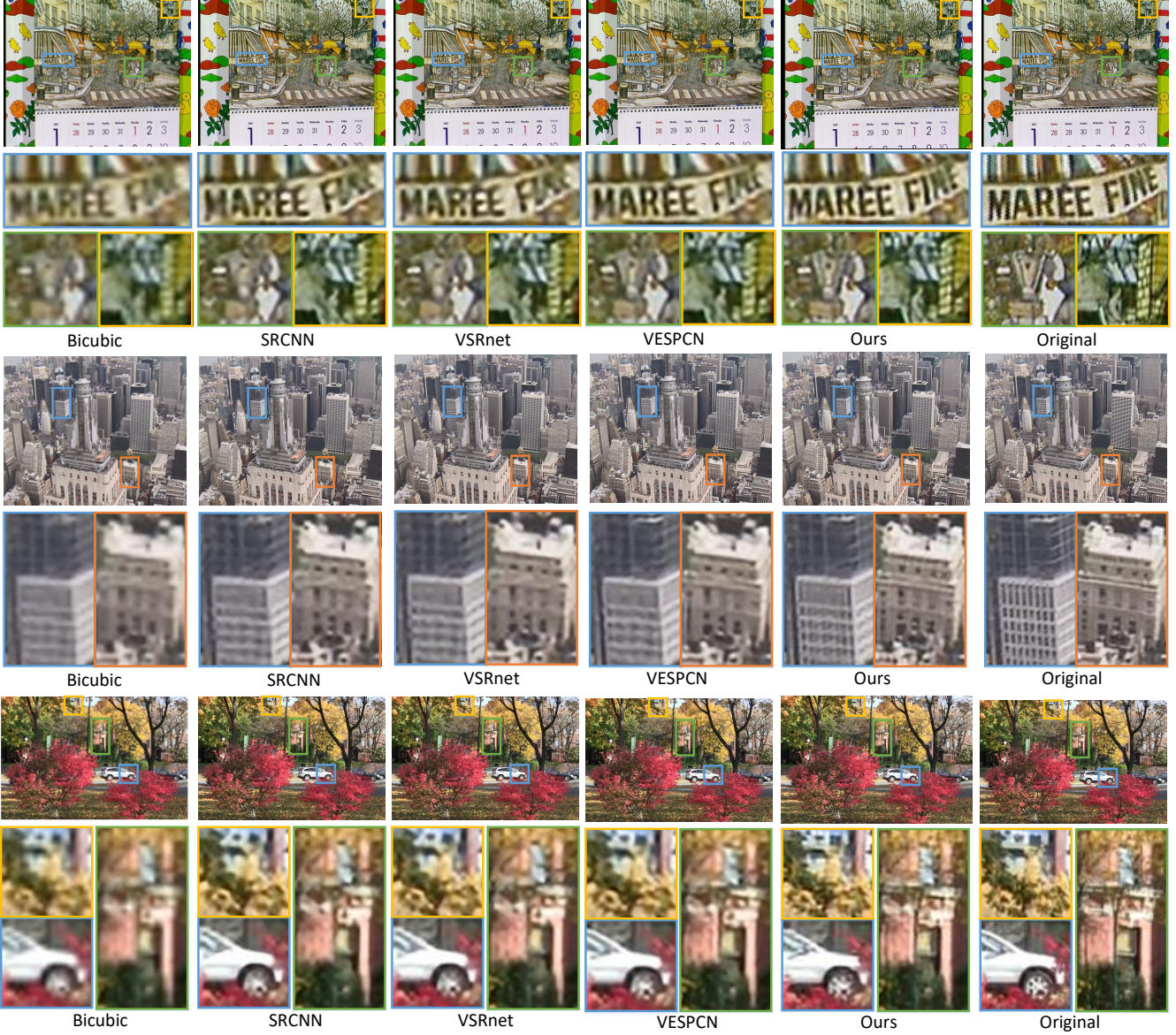| | Image Super-resolution | | | | | | Video Super-resolution | | | | 3DSRnet (**Ours**) | | | |
| | Bicubic | | SRCNN [10] | | VDSR [12] | | VSRnet [3] | | VESPCN [2] | | *smallSet* | | *largeSet* | |
| $\times$ | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 2 | 28.43 | 0.8685 | 30.72 | 0.9176 | 31.44 | 0.9257 | 31.30 | 0.9278 | - | - | 31.98 | 0.9386 | **32.25** | **0.9410** |
| 3 | 25.29 | 0.7341 | 26.54 | 0.7932 | 26.84 | 0.8096 | 26.79 | 0.8098 | 27.25 | 0.8447 | 27.64 | 0.8476 | **27.70** | **0.8498** |
| 4 | 23.79 | 0.6342 | 24.71 | 0.6923 | 24.96 | 0.7121 | 24.84 | 0.7049 | 25.35 | 0.7557 | 25.46 | 0.7498 | **25.71** | **0.7588** |



Figure 7. Comparisons with the state-of-the-art methods for scale factor 3.

shows higher performance with average 0.45 dB and 0.36 dB, respectively for scales 3 and 4 compared to the best performance version (9L-E3-MC) of VESPCN [2] which outperformed its 3D-CNN based video SR version. Furthermore from Table 4, our 3DSRnet performs well on all the four sequences without bias toward certain types of videos. Subjective comparisons for the image and video SR methods in Table 5 are shown in Fig. 7.

Figure 8. Effect of the SF subnet for input with scene change

## 4.4. Scale

Although efficient, a disadvantage of the multi-channel output model [1] is that separate networks have to be trained for different scale factors since the number of output channels should be $scale^2$ - scale to the power of 2. Nevertheless, 3DSRnet with a four-channel output for scale 2 can be trained as a single model for different scales. Specifically, for scale 3, the input frames are first up-scaled by 1.5 times using a bicubic filter and then are fed to the 3DSRnet with scale 2. Similarly, for scale 4, the up-scaled input frames of 2 times are used. For training, we use a dataset that contains a mixture of subimages of all scales 2, 3 and 4 denoted as $sub_2$, $sub_3$ and $sub_4$, respectively, where $sub_3$ and $sub_4$ are up-scaled by 1.5 and 2 times to match the size of $sub_2$. Table 6 shows the PSNR performance of the single model trained for all scales 2, 3 and 4 for the Vidset4 dataset. The single model showed the same PSNR performance compared to the separately trained models for scales 2 and 3, but exhibited a slightly higher performance with average 0.2 dB higher in PSNR for scale 4. The single model benefits from the data of various characteristics having diverse frequency ranges even though it is not devoted to learn the training set of a certain scale.

Table 6. Experiment on scale factors. On the 2nd row are the networks trained separately for each of the scale factors. On the 3rd row are the test results of a single network trained for all scales

| Model | ×2 | ×3 | ×4 |
|---|---|---|---|
| Separate | 31.82 | 27.43 | 25.27 |
| Single | 31.82 | 27.43 | 25.47 |

*PSNR(dB) performance

## 4.5. Scene Change Detection and Frame Replacement Subnet

Scene change often occurs in video sequences, but little attention has been given in video SR. Without a proper treatment for scene change in input frames, performance degradation is inevitable due to the presence of irrelevant frames. Therefore, in the case of scene change, we swap the unrelated frames with the temporally closest frames of the same scene using the SF subnet introduced in Section 3.2. It improves the quality of the output frames significantly. Fig. 8 shows the qualitative and quantitave results of our 3DSRnet with and without frame replacement for input with scene change. As seen in Fig. 8 (c), if the disparate frames are replaced with zeros, the performance tends to severely drop.

Let $F_n$ a series of frames in a sliding time window where the $n$-th frame $f_n$ is the first frame just after scene change. As illustrated Fig. 8 (c), the current frames of $F_1$, $F_2$ and $F_3$ correspond to the middle red frames, and those of $F_4$, $F_5$ and $F_6$ to the yellow frames. $f_1$ of $F_2$, and $f_1$ and $f_2$ of $F_3$ are replaced with $f_2$ of $F_2$, and $f_3$ of $F_3$, respectively. Similarly, $f_4$ and $f_5$ of $F_4$, and $f_5$ of $F_5$ are replaced with $f_3$ of $F_4$, and $f_4$ of $F_5$, respectively. Note that $F_1$ and $F_6$ do not contain any scene change. As shown in Fig. 8 (c), the PSNR starts to drop from $F_1$ to $F_3$, and from $F_6$ to $F_4$. When the SF subnet is incorporated, the PSNR values without the SF subnet are enhanced by average 0.39, 0.46, 0.32, and 0.25 dB for $F_2$, $F_3$, $F_4$ and $F_5$, respectively. Table 7 shows the detection accuracy of scene change by the SF subnet architecture with two and three layers. Even with the

Table 7. Scene change detection accuracy of 2-layer and 3-layer networks. The test set contains 5,240 samples evenly belonging to the five classes

| Network | 2-Layer | 3-Layer |
|---|---|---|
| Detection Accuracy (%) | 99.886 | 99.905 |

9

two-layer SF subnet, the detection accuracy of 99.89% was obtained.

## 4.6. Inference Time

The inference time on an NVIDIA TITAN X GPU is 166 ms and 788 ms for the scale factor of 2 and 4, respectively, to upscale an image of $960 \times 540$ resolution from the input of five frames.

## 5. Conclusion

We propose the 3DSRnet, a video SR method that effectively captures spatio-temporal information of LR input frames in reconstructing HR frames throughout deep 3D convolution layers with temporal depth constantly maintained, all without prior motion alignment. The proposed 3DSRnet employs residual learning with the sub-pixel output structure and prevents severe performance drop due to scene change in the multiple input frames by adopting a simple classification network. The experimental results shows that our proposed 3DSRnet outperformed the most state-of-the-art image and video SR methods by maximum 0.45 dB in PSNR.

## References

[1] Shi, W., Caballero, J., Huszár, F., Totz, J., Aitken, A.P., Bishop, R., Rueckert, D., Wang, Z.: Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2016) 1874–1883

[2] Caballero, J., Ledig, C., Aitken, A., Acosta, A., Totz, J., Wang, Z., Shi, W.: Real-time video super-resolution with spatio-temporal networks and motion compensation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2017)

[3] Kappeler, A., Yoo, S., Dai, Q., Katsaggelos, A.K.: Video super-resolution with convolutional networks. IEEE Transactions on Computational Imaging 2(2) (2016) 109–122

[4] Timofte, R., De, V., Van Gool, L.: Anchored neighborhood regression for fast example-based super-resolution. In: Computer Vision, 2013 IEEE International Conference on, IEEE (2013) 1920–1927

[5] Timofte, R., De Smet, V., Van Gool, L.: A+: Adjusted anchored neighborhood regression for fast super-resolution. In: Asian Conference on Computer Vision, Springer (2014) 111–126

[6] Glasner, D., Bagon, S., Irani, M.: Super-resolution from a single image. In: Computer Vision, 2009 IEEE 12th International Conference on, IEEE (2009) 349–356

[7] Yang, J., Wright, J., Huang, T.S., Ma, Y.: Image super-resolution via sparse representation. IEEE Transactions on Image Processing 19(11) (2010) 2861–2873

[8] Schulter, S., Leistner, C., Bischof, H.: Fast and accurate image upscaling with super-resolution forests. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2015) 3791–3799

[9] Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: Advances in Neural Information Processing Systems. (2012) 1097–1105

[10] Dong, Chao and Loy, Chen Change and He, Kaiming and Tang, Xiaoou.: Learning a deep convolutional network for image super-resolution. In: Proceedings of the European Conference on Computer Vision. (2014) 184–199

[11] Dong, C., Loy, C.C., He, K., Tang, X.: Image super-resolution using deep convolutional networks. IEEE Transactions on Pattern Analysis and Machine Intelligence 38(2) (2016) 295–307

[12] Kim, J., Kwon Lee, J., Mu Lee, K.: Accurate image super-resolution using very deep convolutional networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2016) 1646–1654

[13] Karpathy, A., Toderici, G., Shetty, S., Leung, T., Sukthankar, R., Fei-Fei, L.: Large-scale video classification with convolutional neural networks. In: Proceedings of the IEEE conference on Computer Vision and Pattern Recognition. (2014) 1725–1732

[14] Tran, D., Bourdev, L., Fergus, R., Torresani, L., Paluri, M.: Learning spatiotemporal features with 3d convolutional networks. In: Computer Vision, 2015 IEEE International Conference on, IEEE (2015) 4489–4497

[15] Ji, S., Xu, W., Yang, M., Yu, K.: 3d convolutional neural networks for human action recognition. IEEE Transactions on Pattern Analysis and Machine Intelligence 35(1) (2013) 221–231

[16] Teivas, I.: Video event classification using 3d convolutional neural networks. Master's thesis, Tampere University of Technology (2016)

[17] Glorot, X., Bordes, A., Bengio, Y.: Deep sparse rectifier neural networks. In: Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics. (2011) 315–323

[18] Huang, Z., Wang, L., Meng, G., Pan, C., et al.: Image super-resolution via deep dilated convolutional networks. In: Proceedings of the IEEE International Conference on Image Processing. (2017)

[19] Wang, Q., Fan, H., Cong, Y., Tang, Y.: Large receptive field convolutional neural network for image super-resolution. In: Proceedings of the IEEE International Conference on Image Processing. (2017) 958–962

[20] Shi, W., Jiang, F., Zhao, D.: Single image super-resolution with dilated convolution based multi-scale information learning inception module. In: Proceedings of the IEEE International Conference on Image Processing. (2017)

[21] He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. (2016) 770–778

[22] Bengio, Y., Simard, P., Frasconi, P.: Learning long-term dependencies with gradient descent is difficult. IEEE Transactions on Neural Networks **5**(2) (1994) 157–166

[23] Glorot, X., Bengio, Y.: Understanding the difficulty of training deep feedforward neural networks. In: Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics. (2010) 249–256

[24] Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: Proceedings of the IEEE International Conference on Learning Representations. (2015)

[25] Vedaldi, A., Lenc, K.: Matconvnet: Convolutional neural networks for matlab. In: Proceedings of the 23rd ACM international conference on Multimedia, ACM (2015) 689–692

[26] Liu, C., Sun, D.: On bayesian adaptive video super resolution. IEEE Transactions on Pattern Analysis and Machine Intelligence **36** (2014) 346–360

[27] Liao, R., Tao, X., Li, R., Ma, Z., Jia, J.: Video super-resolution via deep draft-ensemble learning. In: Proceedings of the IEEE International Conference on Computer Vision. (2015) 531–539

[28] Liu, D., Wang, Z., Fan, Y., Liu, X., Wang, Z., Chang, S., Huang, T.: Robust video super-resolution with learned temporal dynamics. In: Proceedings of the IEEE International Conference on Computer Vision. (2017)