

Real-Time Video Scaling Based On Convolution Neural Network Architecture

S Safinaz*¹, AV Ravi kumar²

¹sir.mvit, VTU, Bangalore, India

²Sjbit, VTU, Bangalore, India

*Corresponding author, e-mail: safinaz.mvit@rediffmail.com

Abstract

In recent years, video super resolution techniques becomes mandatory requirements to get high resolution videos. Many super resolution techniques researched but still video super resolution or scaling is a vital challenge. In this paper, we have presented a real-time video scaling based on convolution neural network architecture to eliminate the blurriness in the images and video frames and to provide better reconstruction quality while scaling of large datasets from lower resolution frames to high resolution frames. We compare our outcomes with multiple exiting algorithms. Our extensive results of proposed technique Rem CNN (Reconstruction error minimization Convolution Neural Network) shows that our model outperforms the existing technologies such as bicubic, bilinear, MResNet and provide better reconstructed motioning images and video frames. The experimental results shows that our average PSNR result is 47.80474 considering upscale-2, 41.70209 for upscale-3 and 36.24503 for upscale-4 for Myanmar dataset which is very high in contrast to other existing techniques. This results proves our proposed model real-time video scaling based on convolution neural network architecture's high efficiency and better performance.

Keywords: Image scaling, Convolution Neural Network, Super Resolution

Copyright © 2017 Institute of Advanced Engineering and Science. All rights reserved.

1. Introduction

The In recent years, high definition devices such as HDTV (High-definition television), Smart phones, LAPTOPS, iPad, MacBook Pro and UHD TV (Ultra-high-definition television) have gained immense popularity due to its high resolution quality. Therefore there is an extensive demand of super-resolution in this modern era.

Therefore in recent years Super resolution becomes one of most vital technique for video editing and post-processing applications. Super-resolution is a technique of enhancing the low resolution images or video frames into high resolution frames and images. Super Resolution approach uses neighboring pixels to recover the lost pixels and provide better quality [1]. In many applications such as medical [2], satellite imaging [3], surveillance [4], HDTV [5], video coding or decoding [6-8], stereoscopic video processing [9], [10] and face recognition [11] the use of super resolution becomes mandatory requirement.

Super resolution approach is use to extract high-frequency information from the images and video frames with low resolution quality to reconstruct the original image by eliminating the ringing effect [12]. Hence, Super resolution technique needs high amount of accuracy and speed for the processing of video frame sequences and images. Earlier techniques such as Lanczos, bilinear, and bi-spline provides poor quality of images and video frames with number of visual artifacts like ring, blocking and blurring. However, they are cost efficient and can be easily implemented on chip. As a result of poor resolution they cannot provide required precise high quality of images and video frames. Many issues occurs in the hardware implementation of video scaling such as high computational complexity, large memory requirements, requirement of high resolution quality video, pixel replication and redundancy in pixels. Therefore, this motivates us to implement our video scaling model via software.

Therefore, to eliminate these drawback, in recent years a high resolution CNN (Convolution Neural Network) technique [13, 14] come in the existence. The most dynamic advantage of CNN is that it can easily train with large datasets such as ImageNet [15] and Myanmar dataset [16] by using parallel computing on GPU. These datasets are very bulky in

size which can be a challenging aspect for other existing techniques. CNN techniques are much faster than the conventional techniques due to its easy training and pure feed-forward methods. However, still most of the existing techniques cannot reconstruct the video frames as efficiently as required. Therefore, in this paper, we present a real-time video scaling approach based on *RemCNN* technique to provide high resolution scaling for images and video frames.

Our proposed technique *RemCNN* provides better efficiency and performance in contrast to the existing approaches by eliminating blurriness in the images or video frames to recover the original images and its information. In practical, the key reason of noise occurrence is the difference between the training samples of the training datasets and testing samples of actual application scenes. Therefore, to eliminate this type of noise the proper classification of actual application scenes through Super Resolution (Scaling) approach is necessary so that training samples becomes more similar to that of actual content [17]. In recent years many high resolution devices such as TV (Televisions), laptops and mobile phones developed. However, still many issues such as bulk storage, poor quality and transmission overhead faced by the subscribers. Therefore, to counter these type of conventional issues our proposed *RemCNN* can be prove very vital technique to help researchers and industries considering the current scenarios.

Video Scaling techniques can be partitioned into two parts such as multi-frame and single-frame based approaches [18-19]. Single image based approach mostly utilizes interpolation or example techniques due to their least computational cost. However, in this single-frame based approach resources becomes limited which reduces the system performance hence image quality. Therefore, Video Scaling with multi-frame based approach becomes an import aspect for current scenarios in real time to get better quality reconstructed image. Multi-frame based Video Scaling consists of either reconstruction approach or example based approach or combination of them. However, reconstruction approach provides better fidelity but cannot handle large datasets and large motions. On the other hand, Example-based approaches comes with better performance but mostly depends on quality training [20-21]. Therefore to counter these problems our proposed video scaling approach is highly capable which rely upon *RemCNN* (Reconstruction error minimization Convolution Neural Network) architecture. In this proposed model we use sparse coding reconstruction technique to eliminate the error which generated after feature extraction. We use *SReLU* (Sparse Rectified Linear Unit) to describe non-linearity. Sparse Coding Based Architecture (*SCA*) considered to provide better complex relationship between input low resolution images and its generated output high resolution images

However, previous studies [22-23] consist of some limitations related to its high resolution and image reconstruction when upscaling factor increases. Previous experimental results demonstrate that the efficiency of a system drastically decreases whenever upscaling factor increase. This is due to high frequency component of an image is difficult to extract when scaling factor increases as noise and blurriness level also increases. Therefore our model concentrates on maintaining the efficiency of a system even if upscaling factor further increases. However, to provide better efficiency we apply parallel computing on GPU using CAFFE framework regardless of its upscaling factor. Our experimental results demonstrates that the performance of our video scaling model with *RemCNN* architecture outperforms the existing techniques in terms of scaling factor enhancement outcomes, quality high resolution, noise elimination and precise image reconstruction.

This paper is organized in following sections which are as follows. In section 2, we describe about the video scaling issues and how they can be eliminate by our proposed model. In section 3, we described our proposed methodology. In section 4, experimental results and evaluation shown and section 5 concludes our paper.

2. Video Scaling Issues

Explaining There are many types of issues which can occur while scaling (either upscaling or downscaling) of images and video frames. In [24], Image Fusion and Super-Resolution with Convolutional Neural Network adopted to eliminate blurriness and provide sharp images for digital photography. In this process author Zhong J faces pixel level image fusion issues. In [25], 3D Video Super-Resolution Using Fully Convolutional Neural Networks has been proposed to sort out redundancy, degradation in quality of fused image and huge data size

problems. In [26], Video Super-Resolution with Convolutional Neural Networks adopted to eliminate the problems of video super-resolution. In this paper author faces problem of ill posed in reconstruction of high dimension super resolution image and training of large datasets is also a vital issue. In [18] Image super-resolution: The techniques, Applications, and future provided to review the recent super resolution works and its applications. The biggest challenge face by author Linwei Yue is that to maintain the quality of resolution in motioning conditions. In [27], Real-Time Single Image and Video Super-Resolution Using an Efficient Sub-Pixel Convolutional Neural Network applied. In this paper author experiences global ill-posed reconstruction SR problem and Single Image Super Resolution (SISR) problem which increases the computational complexity of the model. In [28], Image Super-Resolution Based on Convolution Neural Networks Using Multi-Channel Input proposed to get better feature extraction to reconstruct the image. The problems faced by author in achieving this objective are gradient exploding and vanishing and ill posed super resolution problem. In [29], high-quality video/image super-resolution accelerated using GPU to get better performance. In this paper the biggest challenge for the author is the running speed requirement for 4K video processing. In [30], Visualizing and Understanding Convolutional Neural Networks proposed to get better feature extraction and enhancement in image quality. There are two challenges such as training of large datasets like ImageNet dataset with elimination of error and poor capturing of pixels by higher layers are widely faced by the author. In [31], On Bayesian Adaptive Video Super Resolution model presented to get better high resolution reconstructed image with great feature extractions. In this paper, author faces performance degradation issues whenever scaling factor increases. In [32], learning a Mixture of Deep Networks for Single Image Super-Resolution model ill-posed, complex mapping of low-resolution images and inverse image recovery faced by author.

In our paper, the proposed model compared with many existing Super Resolution Video scaling approaches based on CNN framework and there are multiple stages such as shrinking, mapping with sparse coding last layer on which the Convolution Neural Network (CNN) framework rely upon. This stage helps to eliminate the above mentioned poor quality and global ill-posed image reconstruction issues in existing approaches. Dataset videos such as Myanmar video tested with our model and the testing outcomes describes that it can quickly reconstruct the precise information of the video datasets.

The performance of the Super Resolution video scaling architecture significantly increases by using CNN framework. To further improve the performance of the model and speed up the large datasets GPU computing used on a CAFFE framework. CAFFE frameworks not only accelerate the speed of large datasets but also increases the reconstruction quality of images and video frames. The performance of the system remains same in our system regardless of upscaling factor due to fast parallel computing and sparse coding reconstruction architecture. Sparse coding reconstruction technique helps to eliminate sufficient amount noise in image pixels and ill posed problem and reconstruct an efficient original high resolution image.

3. Proposed Methodology

In this section, it is explained the results of research and at the same time is given the comprehensive discussion. Results can be presented in figures, graphs, tables and others that make the reader understand easily [3], [11]. The discussion can be made in several sub-chapters.

3.1. Video Scaling using CNN architecture

In recent years, Convolution neural networks (CNN) gains extreme popularity due to its large success in the field of image or video scaling and image classification [33-34]. CNN can also be easily applied in the fields of face detection [35], pedestrian recognition [36] and object detection [37-38]. CNN provides fast computation for large training database such as ImageNet [15], Myanmar [16] and videoset4 [31]. There are multiple factors which make CNN architecture efficient and help in enhancing the performance of the system.

- a. It helps in the implementation of the training datasets on the efficient and powerful GPU [34] framework such as CAFFE.
- b. It uses ReLU (Rectifier Linear Unit) [36] to provide better performance and fastening speed in training and testing of datasets.
- c. It can easily train large datasets like Myanmar datasets [16].

3.2. Image Reconstruction Architecture

In recent years, precise image reconstruction from low-level resolution to high-level resolution image becomes a mandatory requirement. In previous work many techniques or approaches are applied to reconstruct a better quality image. However very few techniques are able to provide required high resolution reconstructed image. One technique, which shown high accuracy outcomes and better PSNR performance for image reconstruction, is *RemCNN* (Reconstruction error minimization Convolution Neural Networks). In this paper, to compute large training datasets with ultra-high speed, GPU computing used in CAFFE framework. To make our system more precise and eliminate sufficient amount of noise from the image or video frame we apply here sparse coding reconstruction technique for a CNN architecture.

The architectural viewpoint for sparse coding reconstruction method is given in Figure 1 which shows the architecture diagram of reconstruction of image. Consider a single low-dimension video frame. In our proposed model patch based feature extracted for each frame in a video. Then all the frames are down-sampled to the intermediate frames. Then for each frame mapping is require. Then frames are up sampled to the desired size. The difference of up sample and down sample frames fed to sparse coding image reconstruction block to reconstruct image to the original quality. Our proposed model outperforms existing techniques by eliminating the error present in the up sampled image and down sampled image. The reconstruction of image or video frames consists of total five stages in our proposed model.

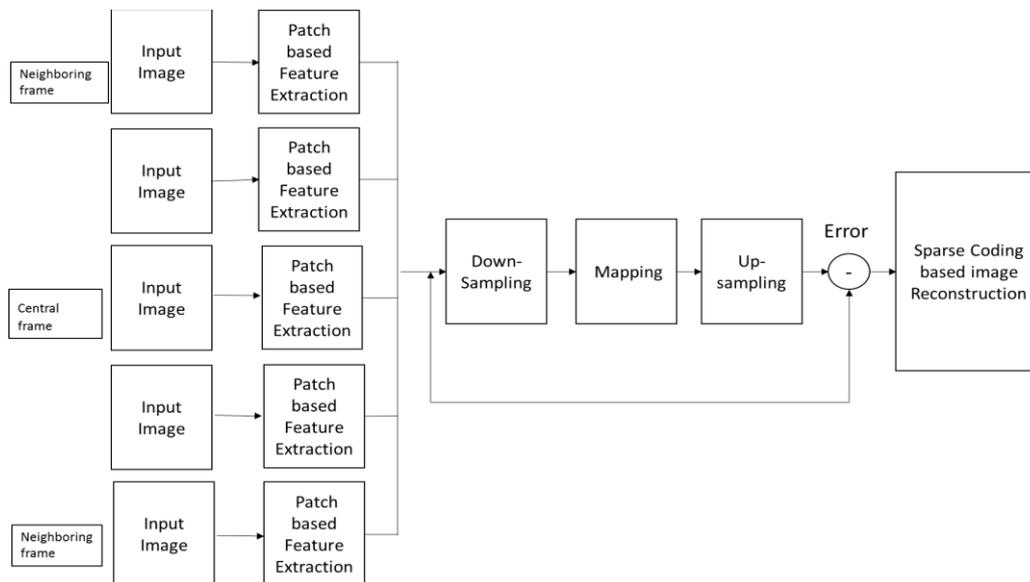


Figure 1. Sparse coding based image reconstruction architecture

There are multiple stages in our proposed model image reconstruction architecture which are as follows.

Patch based feature Extraction: Our proposed *RemCNN* technique first performs patch based feature extraction on each original video frame without interpolation. We represents our input image as X_s . Our input image X_s convoluted to a group of filters to get high dimensional feature vector for each frame. In our model group of filters consists of multiple parameters such as S_1, F_1, M_1 . As our model perform feature extraction directly on the original frames, the filter size of first sheet S_1 can be as $S_1 = 5$. The number of channels we have adopted here out of *YCbCr* is only channel Y hence the number of channel $M_1 = 1$. Here, F_1 is the number of feature dimension which has to determine. This feature dimension of first layer can be presented as $Conv(5, s, 1)$. Here s can be represented as the first sensitive variable.

Down-Sampling: In existing techniques after feature extraction directly mapping presented. Then high dimension features converted into the high resolution features. This

increases computational overhead of the system and degrades the performance due to large size of s .

Therefore to eliminate this drawback of the existing techniques here we first down sample the features extracted by the video frames. This approach can also be observed in high-level vision methods to decrease computational cost.

For the same concern we have down sampled the features of all the layers to decrease the feature dimension s . The filter size for second layer considered as $S_2 = 1$ to perform linearly with features. The feature dimension for second layer can be presented as $F_2 = g \ll s$. Now feature dimensions are decreased from s to g . Here, g is the second sensitive variable which calculates the amount of downsampling. This feature dimension of second layer (1×1) can be presented as $Conv(1, g, s)$. This technique reduces large amount feature dimension.

Mapping: It is the most vital phase of this proposed algorithm which enhances the performance of the model. This is a non-linear type mapping. In mapping, width and depth are two factors which are most affected. Here, width represent the number of filters present in a layer and depth represents the total number of layers. This operation perform non-linear mapping on each high-dimensional feature. In existing techniques mapping experiments not implemented on large deep networks which helps us to create a more significant non-linear mapping layer. To achieve this we consider a medium filter of size $S_3 = 3$. Then, to provide better efficiency we utilize multiple 3×3 layers. The complexity and accuracy performance calculated by a sensitive variable d . Each mapping layer consists of similar number of filters $F_3 = g$. This non-linear mapping can be presented as $Conv(3, g, g)$.

Up-sampling: It is the reverse procedure of the down-sampling. To decrease the feature dimensions down-sampling used which helps in the reduction of computational complexities and produces a high quality video frame or image. Therefore, to generate a high quality image after mapping an up-sampling layer introduced. To retain synchronization between both the layers down-sampling and up-sampling we implemented 1×1 layers. As it is an inverse of down-sampling, the up-sampling layer can represented as $Conv(1, s, g)$. This layer increases the performance of the system.

Sparse coding based image reconstruction: The final part of the image reconstruction is sparse coding based image reconstruction which used to reform a high quality image by eliminating the error produced in up-sampling and down-sampling. Then the outcome (weight parameter) is directly a reformed image with high quality. Here we have taken 9×9 filter layers and the sparse coding layer can be presented as $SparseCode(9, 1, s)$.

3.3. Sparse Rectified Linear Unit (*SReLU*)

After each layer, *ReLU* (Rectified Linear Unit) used for the activation function. In our model we have used Sparse Rectified Linear Unit (*SReLU*) instead of conventional *ReLU*. The activation function for *SReLU* can be

$$f(y_j) = \max(y_j, 0) + b_j \min(0, y_j) \quad (1)$$

Here y_j is the input for the activation function f , j represents the channel and b_j represents the coefficient of negative phase. In existing techniques b_j kept as zero but in for *SReLU* technique b_j is user-defined. *SReLU* is a key to eliminate the dead features [40] generates in *ReLU* by zero gradient vectors. This helps to test parameters of multiple networks for different designs to its full capacity. Our experimental outcomes demonstrate that the *SReLU* networks is comparatively more efficient and stable. This method increases accuracy and speed as well.

3.4. Modelling to Reduce Computational Complexity and Cost Function

3.4.1. Computational Complexity

In existing techniques the computational complexity remains very high which degrades the overall performance of the system. The reason for high computational complexity and cost function is the use of conventional *ReLU* and drawback in the design architecture. In existing approaches computational complexity can be calculated as:

$$O\{(S_1^2 F_1 + F_1 S_2^2 F_2 + F_2 S_3^2) S_{hr}\} \quad (2)$$

Our proposed model consists of very low computational complexity. This is due to its efficient and accurate modern design architecture and use of sparse rectified linear unit (*SReLU*) which helps in increasing the speed and avoiding time lapse by eliminating the dead features. In our proposed model computational complexity calculated as:

$$O\{(25s + gs + 9dg^2 + sg + 81s)S_{tr}\} = O\{(9dg^2)\} \tag{3}$$

3.4.2. Cost Function

In our model cost function described in terms of MSE (Mean Square Root) function. The following equation represent the cost function which used in previous techniques:

$$\min_{\phi} \frac{1}{F} \sum_{k=1}^F \| F(A_g^k; \phi) - B^k \|_2^2 \tag{4}$$

Here, A_g^k and B^k are the k^{th} low and high resolution image pair in training. ϕ is the parameter of the output system function $F(A_g^k; \phi)$. The efficiency of these parameters are maintained by utilizing standard back propagation approach with stochastic gradient.

3.4.3. Sparse Coding Reconstruction

Network Architecture: Sparse Coding Based Architecture (*SCA*) considered to provide better complex relationship between input low resolution images and its generated output high resolution images. This architecture provides better performance and increases high amount of accuracy. This *SCA* (Sparse Coding Based Architecture) implemented in corporation with neural networks to reconstruct a high resolution image from the original low-resolution image using LIST (Learned Iterative Shrinkage and Thresholding) approach [40]. Figure 2 shows the architectural diagram of Sparse Coding Based Architecture (*SCA*).

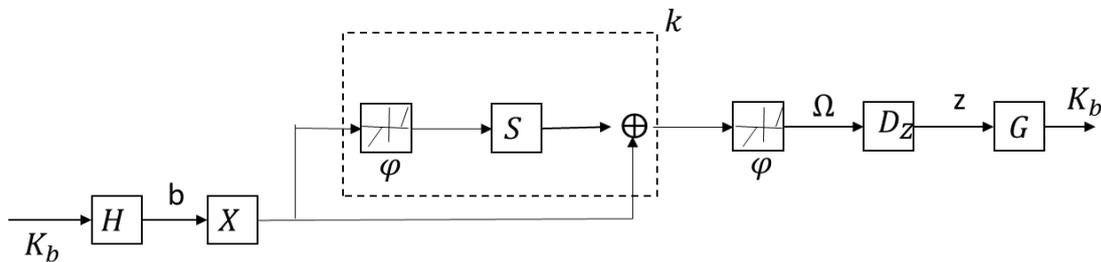


Figure 2. architectural diagram of Sparse Coding Based Architecture (*SCA*)

The objective of our model is to provide high amount of selectivity for the output high resolution image frame by applying *SCA* based LIST approach on each input frame in coordination with neural networks. We use *SReLU* (Sparse Rectified Linear Unit) to describe non-linearity. The use of Reconstruction error minimizer Convolution Neural Networks (*RemCNN*) with *SCA* reduces high amount of computational cost and enhances efficiency with a large extent. For each output high resolution frame a weight map generated based on pixels. Each generated weight-map multiplied with its equivalent pixels for every output frame. Then all the products of frames are summed up to reconstruct an original image frame. The reconstructed original image $F(\mathbf{b}; \Theta)$ can be represented as:

$$F(\mathbf{b}; \Theta) = \sum_k^M X_k(\mathbf{b}; \phi_{\Omega}) \odot F_{c_k}(\mathbf{b}; \phi_{c_k}) \tag{5}$$

Here, \mathbf{b} represent the input low-resolution image, function $X_k(\mathbf{b}; \phi_{\Omega})$ denotes the behavior of generated weight maps and another function $F_{c_k}(\mathbf{b}; \phi_{c_k})$ denotes the output high resolution

frame C_k . Equation (5) represents point wise multiplication of weighted map with its pixels for each reconstructed output image.

Equation (6) the loss elimination between the input low-resolution frame and the estimated output frame in training,

$$\min_{\Theta} \sum_l \|F(\mathbf{b}_l; \Theta) - \mathbf{z}_l\|_2^2 \quad (6)$$

Here, function $F(\mathbf{b}_l; \Theta)$ denotes our output, \mathbf{b}_l is the l^{th} high resolution frame and \mathbf{z}_l represents the corresponding low-resolution image. Θ denotes a group of all parameters of our model. The combination of Equation (5) and (6) provides the cost function of our proposed model.

$$\min_{\varphi_{\Omega}, \{\varphi_{C_k}\}_{k=1}^M} \sum_l \left\| \sum_k X_k(\mathbf{b}_l; \varphi_{\Omega}) \odot F_{C_k}(\mathbf{b}_l; \varphi_{C_k}) - \mathbf{z}_l \right\|_2^2 \quad (7)$$

4. Results and Analysis

We compute our outcomes with the similar dataset (Myanmar) as used in [16] to compare the performance and efficiency of our model to the existing techniques discussed in the related work. Our model is trained on different large dataset like Myanmar [16]. Testing results shows that our model outperforms most of the existing techniques in terms of PSNR and reconstruction efficiency. We have tested our model for different up scaling factors (2, 3 and 4). Our result shows accuracy and reconstruction efficiency increment to a large extent. Our model needs less amount of execution time to provide effective video scaling. Our model implemented on 64-bit windows 10 OS with 16 GB RAM which consists on INTEL (R) core (TM) i5-4460 processor. It consists of 3.20 GHz CPU. We have compared our model with Enhancer [41], Draft-CNN [42], Bayesian [31] and Bayesian-MB [43] and many other existing techniques.

4.1. Implementation Details

We have implemented our extensive experiments on large 4K video Myanmar dataset. In modern era, the availability of 4K monitors is highly increased. Therefore, there is a huge demand of low resolution videos to high-resolution videos in market. These high resolution video can be achieve through upscaling factor. Therefore we have used different upscaling factors to achieve these objectives by measuring performance and accuracy of the model for upscaling factor 2, 3 and 4. Myanmar dataset contains total 57 scenes. In these dataset, 50 scenes used for training and 7 scenes for testing for different up-scaling factors. All the experiments are undertaken on the MATLAB 16b framework in configuration with CAFFE.

4.2. Comparative Study

Here, we have taken 7 scenes for testing out of 57 total scenes in Myanmar dataset video considering upscale-2 as used in [44]. All the scenes are compared to nine most popular existing techniques. Scene-2 represent the famous Myanmar temple which consists of total 594 frames. Our proposed technique *RemCNN* gives 48.0492 dB PSNR. Scene-8 represents Myanmar golden temple which consists of 354 frames. Our proposed technique *RemCNN* gives 36.99 dB PSNR for scene-8 which is little less compare to other existing technique. It is an exceptional case in our proposed model. Scene 18 and 33 represents snake and Buddha temple scenes in Myanmar video. Both the scenes consists of 632 frames and for both scenes our proposed technique *RemCNN* gives highest PSNR as 52.274 and 53.198 dB. Scene 25 and 45 represents yoga scene by a man and horse scenes in Myanmar video. Both the scenes consists of 594 frames and for both scenes our proposed technique *RemCNN* gives PSNR as 47.031 and 49.671 dB. Scene 48 represent tiger scene in Myanmar video. Our proposed technique *RemCNN* gives PSNR as 47.42 dB for this scene. Similarly, same scenes are used to compute PSNR considering upscale-3 and upscale-4. Our proposed technique shows highest PSNR for scene-1 (temple) which consists of 816 frames. The PSNR results for different upscaling factor 2, 3 and 4 are 54.07, 48.96 and 45.05 dB which is much better than the existing techniques. The percentage improvement of our proposed model in contrast to other

conventional techniques is very high. Scene-48 gives highest improvement of 22.17% considering upscale-2. Similarly, scene-45, 33 and 18 gives improvement of 10.21%, 16.37% and 15.22% respectively. However, scene-8 and 25 gives little less accuracy considering upscale-2. Our model gives best improvement result for scene-1 as 28.83% considering upscale-2.

Similarly, Scene-18 gives highest improvement of 15.22% considering upscale-4. Similarly, scene-2, 8, 45 and 48 gives improvement of 3.95%, 12.46%, 2.18% and 14.07 % respectively. However, scene-25 and 33 gives little less accuracy considering upscale-4. Our model gives best improvement result for scene-1 as 28.17% considering upscale-4. Similarly, our model gives best improvement result for scene-1 as 28.15% considering upscale-3. Average PSNR improvement considering upscale-2 and upscale-4 is 7.79% and 4.45%.

Table 1 shows the comparison of different scenes of a Myanmar dataset for multiple existing techniques. The following results shows that our average PSNR result is 47.80474 dB considering upscale-2, 41.70209 dB for upscale-3 and 38.24503 dB for upscale-4 (table 4.3) considering all seven testing scenes which is much better than the existing techniques for MYANMAR dataset. Similarly, Table 2 and 3 represent SSIM (structural similarity index) comparison with recent existing techniques considering upscale-2 and 4 for scene 2, 8, 18, 25, 33, 45 and 48. Table 4 represent SSIM comparison with *VSRnet* and *MCResNet* considering upscale 2, 3 and 4 for scene-1 which is better than existing techniques.

Table 1. PSNR values (in DB) of the SR frame for various methods and test scenes (best results are shown in bold) considering upscale-2

Scenes	bicubic	bi-level	SDMF-B	SDMF-R	MDSF	MDMF-B	MDMF-R	MDMF-B-VT	MDMF-R-VT	Our Proposed
Scene-2	45.27	46.12	46.81	46.41	46.79	47.66	46.86	48.14	47.41	48.0492
Scene-8	38.18	39.94	40.08	40.32	40.34	40.59	40.60	40.98	41.05	36.99
Scene-18	41.43	43.04	43.41	43.69	43.37	43.92	44.19	44.32	44.46	52.274
Scene-25	44.40	46.69	47.52	47.68	47.37	48.45	47.83	49.19	48.59	47.031
Scene-33	40.22	42.95	43.08	43.55	43.27	43.68	44.05	44.49	44.48	53.198
Scene-45	42.43	43.72	44.07	44.18	44.05	44.49	44.28	44.60	44.62	49.671
Scene-48	33.90	36.10	36.20	35.66	36.55	36.67	36.07	36.91	36.64	47.42

Table 2. PSNR values (in DB) of the SR frame for various methods and test scenes (best results are shown in bold) considering upscale-4

Scenes	Bicubic	bi-level [45]	NE+NN LS [46]	NE+LLE [47]	ANR [48]	SR-CNN [49]	Enhance r [41]	Bayesian [31]	MDMF-B-VT [44]	MDMF-R-VT [44]	Our Proposed
Scene-2	39.58	40.50	41.32	41.12	41.32	43.17	40.62	39.18	43.48	42.90	45.2676
Scene-8	32.13	32.46	33.00	32.95	32.81	33.40	32.09	31.73	33.48	33.42	38.2453
Scene-18	35.65	36.37	36.76	36.82	36.76	37.50	36.44	35.70	37.68	37.65	44.44899
Scene-25	36.10	37.02	37.90	37.78	37.49	38.35	37.44	35.34	39.03	38.75	38.96423
Scene-33	32.15	33.44	33.79	33.94	34.00	34.57	34.67	32.14	34.92	34.86	28.45522
Scene-45	36.13	36.71	37.12	37.27	37.35	37.90	37.15	35.76	38.42	38.10	38.87859
Scene-48	27.25	28.03	28.04	28.20	28.26	28.73	27.75	26.76	28.75	28.49	33.45522
Average	34.14	34.93	35.42	35.44	35.43	36.23	35.17	33.80	36.54	36.31	38.24503

Table 3. PSNR values (in DB) of the SR frame for various methods and test scenes Considering Upscale-2,3,4 for MYANMAR dataset

Scenes	Our proposed (upsampling-2)	Our proposed (upsampling-3)	Our proposed (upsampling-4)
Scene-2	48.0492	49.34129143	45.26765851
Scene-8	36.99	41.61521459	38.24530073
Scene-18	52.274	48.24806934	44.44898805
Scene-25	47.031	42.87746065	38.96423093
Scene-33	53.198	31.56951803	28.45521909
Scene-45	49.671	42.19030456	38.87859171
Scene-48	47.42	36.07273785	33.45522214
AVERAGE	47.80474	41.70209	38.24503

Table 4. PSNR comparison for upscale 2, 3 and 4 with MD MFB-VT and MD MFR-VT for scene-1

	MD MFB-VT	MD MFR-VT	Our <i>RemCNN</i>
UPSCALE-2	38.48	40.04	54.07
UPSCALE-3	34.42	35.18	48.96
UPSCALE-4	31.85	32.36	45.05

Table 5. SSIM values (in DB) of the SR frame for various methods and test scenes (best results are shown in bold) considering upscale-2

Scenes	Bicubic	bi-level [45]	NE + NNLS [46]	NE + LLE [47]	ANR [48]	SR- CNN [49]	Enhan cer [41]	Bayesi an [31]	MDMF- B-VT [44]	MDMF- R-VT [44]	Our Proposed
Scene-2	0.9830	0.9879	0.9851	0.9834	0.9857	0.9859	0.9854	0.9874	0.9882	0.9882	0.9971
Scene-8	0.9738	0.9842	0.9824	0.9817	0.9832	0.9852	0.9823	0.9828	0.9884	0.9882	0.9926
Scene-18	0.9738	0.9849	0.9820	0.9816	0.9833	0.9844	0.9844	0.9842	0.9877	0.9884	0.9963
Scene-25	0.9917	0.9961	0.9938	0.9936	0.9952	0.9955	0.9938	0.9954	0.9970	0.9967	0.9960
Scene-33	0.9786	0.9904	0.9879	0.9889	0.9902	0.9907	0.9908	0.9900	0.9937	0.9938	0.9801
Scene-45	0.9718	0.9810	0.9772	0.9776	0.9791	0.9797	0.9764	0.9790	0.9812	0.9823	0.9942
Scene-48	0.9668	0.9808	0.9774	0.9785	0.9799	0.9826	0.9751	0.9770	0.9846	0.9821	0.9872
Average	0.9771	0.9865	0.9837	0.9836	0.9851	0.9863	0.9840	0.9851	0.9887	0.9885	0.9919

Table 6. SSIM values (in DB) of the SR frame for various methods and test scenes (best results are shown in bold) considering upscale-4

Scenes	Bicubic	bi-level [45]	NE+NN LS [46]	NE+LL E [47]	ANR [48]	SR- CNN [49]	Enhan cer [41]	Bayesi an [31]	MDMF- B-VT [44]	MDMF- R-VT [44]	Our Proposed
Scene-2	0.9648	0.9662	0.9691	0.9675	0.9691	0.9703	0.9695	0.9660	0.9737	0.9740	0.9900
Scene-8	0.9013	0.9099	0.9145	0.9187	0.9107	0.9198	0.9121	0.8972	0.9266	0.9250	0.9515
Scene-18	0.9122	0.9209	0.9243	0.9249	0.9243	0.9280	0.9308	0.9183	0.9331	0.9341	0.9837
Scene-25	0.9515	0.9546	0.9622	0.9607	0.9587	0.9633	0.9621	0.9473	0.9702	0.9687	0.9741
Scene-33	0.8899	0.9140	0.9157	0.9188	0.9206	0.9230	0.9304	0.8945	0.9363	0.9374	0.8795
Scene-45	0.9101	0.9155	0.9193	0.9211	0.9226	0.9253	0.9267	0.9083	0.9340	0.9316	0.9617
Scene-48	0.8514	0.8730	0.8710	0.8757	0.8780	0.8883	0.8679	0.8393	0.8921	0.8842	0.9224
Average	0.9116	0.9220	0.9252	0.9268	0.9263	0.9311	0.9285	0.9101	0.9380	0.9364	0.9519

Table 7. SSIM comparison for upscale 2, 3 and 4 with MD MFB-VT and MD MFR-VT for scene-1

	<i>VSRnet</i>	<i>MResNET</i>	Our <i>RemCNN</i>
UPSCALE-2	0.9679	0.9777	0.9973
UPSCALE-3	0.9247	0.9387	0.9956
UPSCALE-4	0.8834	0.8987	99.08

4.3. Image Reconstruction Comparison

Here, we have demonstrated 350th frame of scene-8 as used in all the other existing techniques. The original Myanmar video dataset contains total 57 scenes and its original resolution is 3840 × 2160. We have shown PSNR and image reconstruction quality comparison with all the conventional techniques. The PSNR result (41.615 dB) outperforms all the existing state-of-the-art techniques. From our experimental results it is clearly visible that our reconstruct frame has better reconstruction quality than any other recent existing techniques. Table 8 comparison with different existing techniques.



GROUND TRUTH



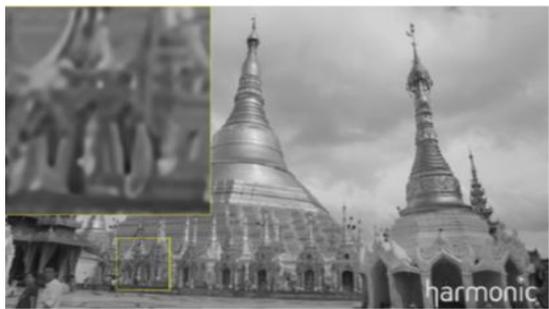
BICUBIC/29.42 dB



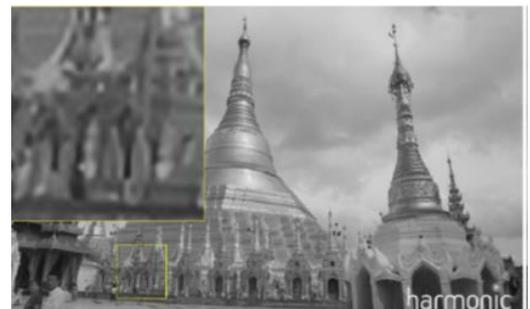
SRCNN[48]/30.77 dB



VDSR[52]/31.07 dB



DRCN[51]/31.05 dB



VSRnet AMC [16]/31.02 dB



VSRnet MC [16]/31.12 dB



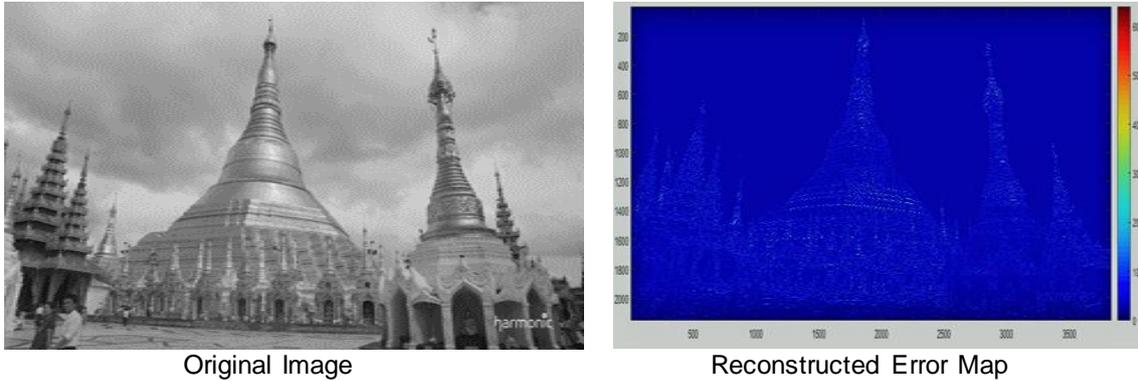
MCRResNet [50]/32.16 dB



ProposedRemCNN/41.615 dB

4.4. Reconstruction Error Map

Here, we reconstruct error map for 350th frame of scene-8 shown in Table 8. Reconstruction map is the combination of planning motion and local motion of background and foreground respectively. The reconstruction error depends on the iterations as the iteration increases, the error become decreases. In the final outcomes error become disappears or become negligible using our proposed technique *RemCNN*. Table 9 reconstructed error map from original image.



4.5. Graphical Analysis

The following graphs shows the comparison between our proposed model and existing approaches MD MFB-VT and MD MFR-VT for upscale 2, 3 and 4 considering Myanmar dataset. Figure 1 shows PSNR comparison considering upscale -2 for the scenes 2, 8, 18, 25, 33, 45 and 48. Figure 2 shows PSNR comparison considering upscale -4 for the scenes 2, 8, 18, 25, 33, 45 and 48. Figure 3 demonstrates PSNR comparison considering upscale-2,3 and 4 for scene-1 with both recent existing techniques MD MFB-VT and MD MFR-VT. PSNR for upscale-2 using our proposed *RemCNN* technique is 54.07 dB, with upscale-3 is 48.96 dB and for upscale-4 is 45.05 db which is very high compare to other techniques.

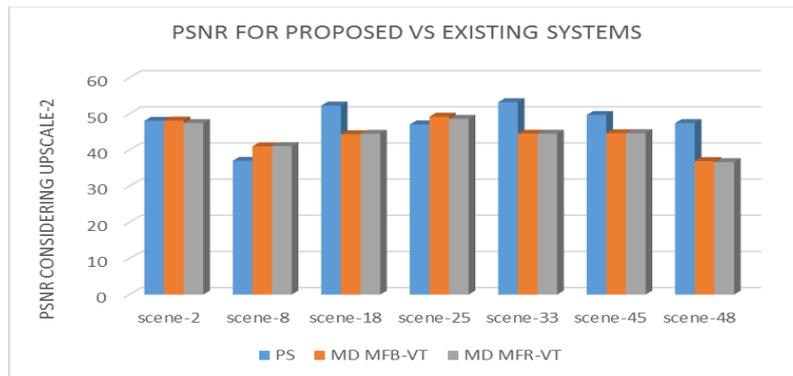


Figure 3. PSNR comparison for proposed vs existing techniques for upscaling factor -2 for Myanmar DATASET

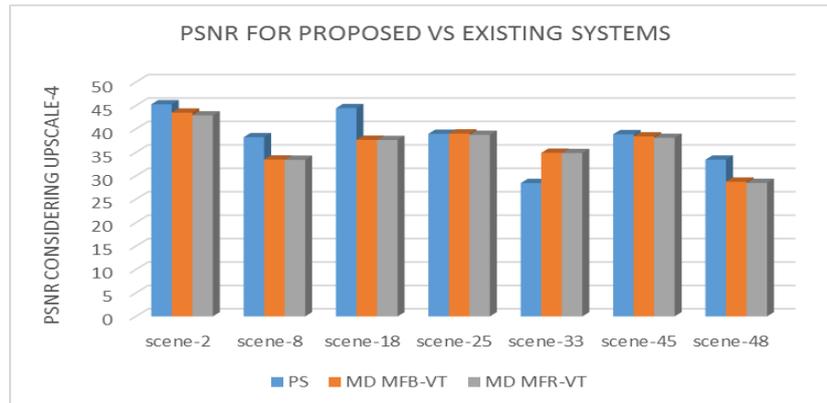


Figure 4. PSNR comparison for proposed vs existing techniques for upscaling factor -4 For Myanmar dataset

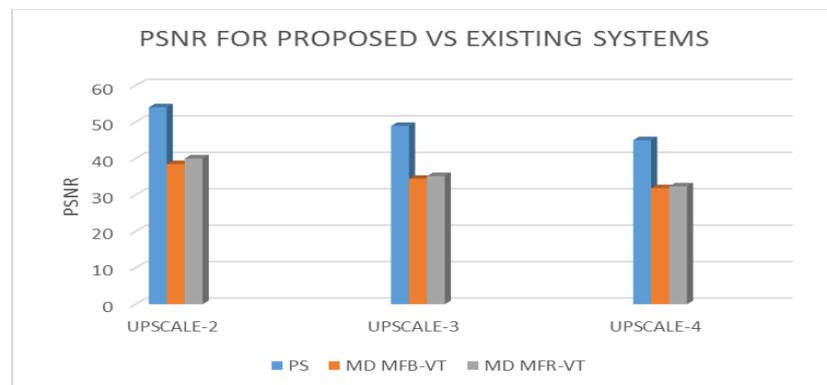


Figure 5. PSNR for proposed vs existing considering upscale factor 2, 3 and 4 for Myanmar dataset

4. Conclusion

In current era, huge demand and popularity of high resolution videos made researchers to carry out work in video scaling field to offer ease of accessibility of high-resolution videos to the subscribers. Therefore, we have introduced a real-time video scaling based on convolution neural network architecture to eliminate the blurriness in the images and video frames and to provide better reconstruction quality while scaling of large datasets. CNN architecture helps us to restore high frequency components of the video frames. Our proposed model can easily train the bulky datasets such as Myanmar and Videose4. Our experimental results shows that our model outperforms many existing techniques in terms of PSNR, fidelity and reconstruction quality. The experimental results shows that our average PSNR result is 47.80474 considering upscale-2, 41.70209 for upscale-3 and 36.24503 for upscale-4 for Myanmar dataset which is very high in contrast to other existing techniques. Our model gives best improvement result for scene-1 as 28.83% considering upscale-2, 28.17% considering upscale-4, 28.15% for upscale-3.

This results proves our proposed model real-time video scaling based on convolution neural network architecture's high efficiency and better performance. Our proposed model can be effectively used in the applications such as medical, satellite imaging, surveillance, HDTV, video coding or decoding, stereoscopic video processing, and face recognition for future purpose to reconstruct efficient images or video frames.

References

- [1] Dong W, Zhang L, Shi G, et al. Image deblurring and super-resolution by adaptive sparse domain selection and adaptive regularization. *IEEE Transactions on Image Processing*. 2011; 20(7): 1838-1857.
- [2] W Shi, J Caballero, C Ledig, X Zhuang, W Bai, K Bhatia, A Marvao, T Dawes, D Oregan and D Rueckert. Cardiac image super-resolution with global correspondence using multi-atlas patchmatch. In K Mori, I Sakuma, Y Sato, C Barillot and N Navab. *Editors, Medical Image Computing and Computer Assisted Intervention (MICCAI)*. 2013; 8151 of LNCS: 9–16.
- [3] MW Thornton, PM Atkinson and DA Holland. Sub-pixel mapping of rural land cover objects from fine spatial resolution satellite sensor imagery using super-resolution pixel-swapping. *International Journal of Remote Sensing*. 2006; 27(3): 473–491.
- [4] Yong Chen, Feng Shuai. Real-time Colorized Video Images Optimization Method in Scotopic Vision. *TELKOMNIKA Indonesian Journal of Electrical Engineering*. 2015; 15(2): 321 ~ 330. DOI: 10.11591/telkomnika.v15i2.8125
- [5] T Goto, T Fukuoka, F Nagashima, S Hirano and M Sakurai. *Super-resolution System for 4K-HDTV*. 2014 22nd International Conference on Pattern Recognition. 2014: 4453–4458.
- [6] BC Song, SC Jeong and Y Choi. "Video super-resolution algorithm using bi-directional overlapped block motion compensation and on-the-fly dictionary training". *IEEE Trans. Circuits Syst. Video Technol.* 2011; 21(3): 274–285.
- [7] EM Hung, RL De Queiroz, F Brandi, KF De Oliveira and D Mukherjee. "Video super-resolution using codebooks derived from key-frames". *IEEE Trans. Circuits Syst. Video Technol.* 2012; 22(9): 1321–1331.
- [8] Z Zhang and V Sze. "Fast: Free adaptive super-resolution via transfer for compressed videos". arXiv preprint arXiv: 1603.08968. 2016.
- [9] J Zhang, Y Cao, ZJ Zha, Z Zheng, CW Chen and Z Wang. "A unified scheme for super-resolution and depth estimation from asymmetric stereoscopic video". *IEEE Trans. Circuits Syst. Video Technol.* 2016; 26(3): 479–493.
- [10] Z Jin, T Tillo, C Yao, J Xiao and Y Zhao, "Virtual-view-assisted video super-resolution and enhancement". *IEEE Trans. Circuits Syst. Video Technol.* 2016; 26(3): 467–478.
- [11] BK Gunturk, AU Batur, Y Altunbasak, MH Hayes and RM Mersereau. Eigenface-domain super-resolution for face recognition. *IEEE Transactions on Image Processing*. 2003; 12(5): 597–606.
- [12] Tappen MF, Russell BC, Freeman WT. *Exploiting the sparse derivative prior for super-resolution and image demosaicing*. In IEEE Workshop on Statistical and Computational Theories of Vision. 2003.
- [13] A Krizhevsky, I Sutskever, and GE Hinton. "Imagenet classification with deep convolutional neural networks". In *NIPS*. 2012: 1097–1105.
- [14] C Szegedy, W Liu, Y Jia, P Sermanet, S Reed, D Anguelov, D Erhan, V Vanhoucke, and A Rabinovich. "Going deeper with convolutions". arXiv preprint: 1409.4842. 2014.
- [15] J Deng, W Dong, R Socher and L Li. "A large-scale hierarchical image database". Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2009: 248–255.
- [16] "Harmonic Inc," <http://www.harmonicinc.com/resources/videos/4kvideo-clip-center>". 2014.
- [17] Zhang L, Zhang L, Mou X, et al. FSIM: a feature similarity index for image quality assessment. *IEEE transactions on Image Processing*. 2011; 20(8): 2378-2386.
- [18] L Yue, H Shen, J Li, Q Yuan, H Zhang, L Zhang. "Image super-resolution: The techniques applications and future". *Signal Process*. 2016; 128: 389-408.
- [19] K Nasrollahi and TB Moeslund. "Super-resolution: a comprehensive survey". *Mach. Vis. & Appl.* 2014; 25(6): 1423–1468.
- [20] S Farsiu, MD Robinson, M Elad and P Milanfar. "Fast and robust multiframe super resolution". *IEEE Trans. Image Process*. 2004; 13(10): 1327–1344.
- [21] M Protter, M Elad, H Takeda and P Milanfar. "Generalizing the nonlocal-means to super-resolution reconstruction". *IEEE Trans. Image Process*. 2009; 18(1): 36–51.
- [22] Pawar Ashwini Dilip, K Rameshbabu. Bilinear Interpolation Image Scaling Processor for VLSI Architecture. *International Journal of Reconfigurable and Embedded Systems (IJRES)*. 2014; 3(3): 104–113. ISSN: 2089-4864
- [23] Z Lin and HY Shum. Fundamental limits of reconstruction based super resolution algorithms under local translation. *IEEE Transaction on Pattern Analysis Machine Intelligence*. 2004; 26: 83 –97.
- [24] Zhong J, Yang B, Li Y, Zhong F, Chen Z. Image Fusion and Super-Resolution with Convolutional Neural Network. In: Tan T, Li X, Chen X, Zhou J, Yang J, Cheng H. (eds) *Pattern Recognition*. CCPR 2016. Communications in Computer and Information Science. Springer, Singapore. 2016; 663.
- [25] Y Xie, J Xiao, T Tillo, Y Wei and Y Zhao. "3D video super-resolution using fully convolutional neural networks". 2016 *IEEE International Conference on Multimedia and Expo (ICME)*, Seattle, WA. 2016: 1-6. doi: 10.1109/ICME.2016.7552931
- [26] A Kappeler, S Yoo, Q Dai and AK Katsaggelos. "Video Super-Resolution with Convolutional Neural Networks". in *IEEE Transactions on Computational Imaging*. 2016; 2(2): 109-122. doi: 10.1109/TCI.2016.2532323

- [27] W Shi et al. "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)*, Las Vegas, NV. 2016: 1874-1883. doi: 10.1109/CVPR.2016.207
- [28] GY Youm, SH Bae and M Kim. "Image super-resolution based on convolution neural networks using multi-channel input". 2016 *IEEE 12th Image, Video, and Multidimensional Signal Processing Workshop (IVMSP)*, Bordeaux. 2016: 1-5. doi: 10.1109/IVMSPW.2016.7528224
- [29] Z Zhao, L Song, R Xie and X. Yang. "GPU accelerated high-quality video/image super-resolution," 2016 *IEEE International Symposium on Broadband Multimedia Systems and Broadcasting (BMSB)*, Nara. 2016: 1-4. doi: 10.1109/BMSB.2016.7521938
- [30] MD Zeiler and R Fergus. Visualizing and understanding convolutional networks. In DJ Fleet, T Pajdla, B Schiele and T Tuytelaars. *Editors, ECCV, Lecture Notes in Computer Science*, Springer, 2014; 8689: 818–833.
- [31] C Liu and D Sun. "On Bayesian Adaptive Video Super Resolution". In *IEEE Transactions on Pattern Analysis and Machine Intelligence*. 2014; 36(2): 346-360. doi: 10.1109/TPAMI.2013.127
- [32] Wang Z, Liu D, Yang J, Han W, Huang T. *Deep networks for image super-resolution with sparse prior*. In: Proceedings of the IEEE International Conference on Computer Vision. 2015: 370-378.
- [33] He K, Zhang X, Ren S, Sun J. Spatial pyramid pooling in deep convolutional networks for visual recognition. In: *European Conference on Computer Vision*. 2014: 346–361.
- [34] Krizhevsky A, Sutskever I, Hinton G. ImageNet classification with deep convolutional neural networks. In: *Advances in Neural Information Processing Systems*. 2012: 1097–1105:
- [35] Sun Y, Chen Y, Wang X, Tang X. Deep learning face representation by joint identification-verification. In: *Advances in Neural Information Processing Systems*. 2014: 1988–1996.
- [36] Szegedy C, Reed S, Erhan D, Anguelov D. Scalable, highquality object detection. arXiv preprint arXiv: 1412.1441. 2014.
- [37] Ouyang W, Luo P, Zeng X, Qiu S, Tian Y, Li H, Yang S, Wang Z, Xiong Y, Qian C, et al. Deepid-net: multi-stage and deformable deep convolutional neural networks for object detection. arXiv preprint arXiv:1409.3505. 2014.
- [38] Ouyang W, Wang X. Joint deep learning for pedestrian detection. In: *IEEE International Conference on Computer Vision*. 2013: 2056–2063.
- [39] Nair V, Hinton GE. Rectified linear units improve restricted Boltzmann machines. In: *International Conference on Machine Learning*. 2010: 807–814.
- [40] Zeiler MD, Fergus R. Visualizing and understanding convolutional networks. In: *ECCV*, Springer. 2014: 818–833.
- [41] Infognition. 2010. Video enhancer. [Online]. Available: <http://www.infognition.com/videoenhancer/>
- [42] R Liao, X Tao, R Li, Z Ma and J Jia. "Video super-resolution via deep draft-ensemble learning". In Proceedings of the IEEE International Conference on Computer Vision. 2015: 531–539.
- [43] Z Ma, R Liao, X Tao, L Xu, J Jia and E Wu. "Handling motion blur in multi-frame super-resolution". in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2015: 5224–5232.
- [44] Q Dai, S Yoo, A Kappeler and AK Katsaggelos. "Sparse Representation-Based Multiple Frame Video Super-Resolution". In *IEEE Transactions on Image Processing*. 2017; 26(2): 765-781. doi: 10.1109/TIP.2016.2631339
- [45] J Yang, Z Wang, Z Lin, X Shu and T Huang. "Bilevel sparse coding for coupled feature spaces," in 2012 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, IEEE. 2012: 2360–2367.
- [46] M Bevilacqua, A Roumy, C Guillemot and ML Alberi-Morel. "Lowcomplexity single-image super-resolution based on nonnegative neighbor embedding". Proceedings of the 23rd British Machine Vision Conference (BMVC). 2012: 135.1–135.10.
- [47] H Chang, DY Yeung and Y Xiong. "Super-resolution through neighbor embedding". In Proceedings of the 2004 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2004. CVPR 2004. IEEE. 2004; 1: I-I.
- [48] R Timofte, V De and L Van Gool. "Anchored neighborhood regression for fast example-based super-resolution". In 2013 *IEEE International Conference on Computer Vision (ICCV)*. IEEE. 2013: 1920–1927.
- [49] C Dong, CC Loy, K He and X Tang. "Learning a deep convolutional network for image super-resolution". In *Computer Vision–ECCV 2014*. Springer. 2014: 184–199.
- [50] D Li, Z Wang. "Video Super-Resolution via Motion Compensation and Deep Residual Learning". in *IEEE Transactions on Computational Imaging*. 99: 1-1. doi: 10.1109/TCI.2017.2671360
- [51] WT Freeman, TR Jones and EC Pasztor. "Example-based superresolution". *IEEE Comput. Graph. and Appl.* 2002; 22(2): 56–65.
- [52] J Kim, JK Lee and KM Lee. "Accurate image super-resolution using very deep convolutional networks". in Proc. IEEE Conf. Comput. Vis. Pattern Recog. 2016: 1646–1654.