

# ENHANCING QUALITY FOR VVC COMPRESSED VIDEOS BY JOINTLY EXPLOITING SPATIAL DETAILS AND TEMPORAL STRUCTURE

Xiandong Meng<sup>1</sup>, Xuan Deng<sup>2</sup>, Shuyuan Zhu<sup>2</sup> and Bing Zeng<sup>2</sup>

<sup>1</sup>The Hong Kong University of Science and Technology

<sup>2</sup>University of Electronic Science and Technology of China

## ABSTRACT

In this paper, we propose a quality enhancement network of versatile video coding (VVC) compressed videos by jointly exploiting spatial details and temporal structure (SDTS). The proposed network consists of a temporal structure fusion subnet and a spatial detail enhancement subnet. The former subnet is used to estimate and compensate the temporal motion across frames, and the latter subnet is used to reduce the compression artifacts and enhance the reconstruction quality of compressed video. Experimental results demonstrate the effectiveness of our SDTS-based method. The code of our proposed method is available at <https://github.com/mengab/SDTS>

**Index Terms**— versatile video coding, spatial-temporal structure, motion compensation, quality enhancement.

## 1. INTRODUCTION

Versatile video coding (VVC) [1] achieves a higher compression performance compared with the High Efficiency Video Coding (HEVC) [2]. Similar to previous video coding standards, the VVC also employs a hybrid scheme which includes the block-based prediction and transform coding to compress videos. Due to the quantization of the transform coefficients in each small block, the artifacts, such as the blocking artifacts and the ringing effects, usually exist in the compressed videos, especially at the low bit-rate. Therefore, it is necessary to enhance the quality of compressed video.

In this work, we focus on the quality enhancement for the compressed video signals based on the latest convolutional neural network (CNN) method. The video quality enhancement may be regarded as the extension of the image quality enhancement in the temporal dimension. Such an extension introduces more prior information which can be used to potentially improve the quality of each individual frame. However, there still exist some challenges to utilize these information to construct an efficient CNN-based solution. First, removing compression artifacts from videos requires the understanding of not only the spatial context of the single frame but also the motion information across frames. Second, although it is possible to find missing content of the same scene

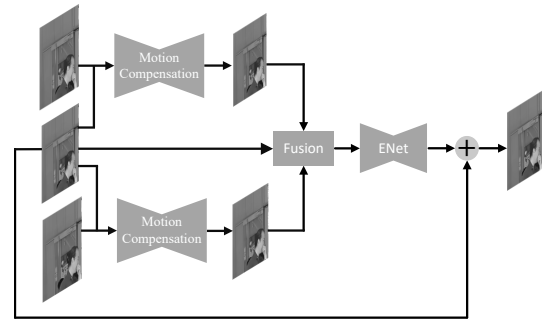


Fig. 1. The framework of our proposed SDTS-based method

or object in adjacent frames, the interference information will be introduced to the target frame if the adjacent frames are directly input to network as a reference. Third, due to the quality fluctuation across compressed video frames, it is very difficult to enhance all video frames with a single model.

We propose a novel end-to-end deep learning architecture in this work to tackle the above issues. The framework of the proposed network is shown in Fig. 1, which consists of a temporal fusion subnet and a spatial detail enhancement subnet. The first subnet is utilized to estimate and compensate the temporal motion across frames, and the second one is employed to reduce the compression artifacts. In addition, as pointed out in [3, 4, 5, 6] that the low-quality frames (LQFs) may be enhanced using the adjacent high-quality frames (HQFs), we also employ the adjacent HQFs as a reference to enhance the low-quality frames. The experimental results demonstrate the performance of the SDTS-based method.

## 2. RELATED WORK

Deep learning has been successfully applied to video super-resolution [7], deblurring [8] and inpainting [9], and can also be employed to enhance the quality of compressed image/video [6, 10, 11, 12, 13, 14, 15, 16, 17]. Particularly, Dong *et al.* [11] firstly proposed ARCNN to reduce the JPEG artifacts of images. Later on, DnCNN [16] and MemNet [12] were proposed for image restoration, including the image quality enhancement. For the quality enhancement of com-

pressed video, VRCNN [10] was proposed as a variable-filter-size residue-learning network [18] for the post-processing of HEVC intra coding. Wang *et al.* [13] developed a Deep CNN-based Auto Decoder (DCAD), which contains 10 CNN layers to reduce the distortion of compressed video. These methods were proposed based only on the prior information of a single frame, so the enhancement performance is still limited. To tackle this problem, Yang *et al.* [6] proposed a MFQE model with multi-frame input for quality enhancement of HEVC compressed video in which the information of neighboring key frames was considered. Meanwhile, Meng *et al.* [15] designed a multi-frame guided attention network by taking advantage of the intra-frame prior information and multi-frame information to enhance the quality of the HEVC compressed video. The experimental results of [6] and [15] have demonstrated that utilizing the multi-frame information to build up the network for video quality enhancement can achieve excellent performance.

### 3. OUR PROPOSED METHOD

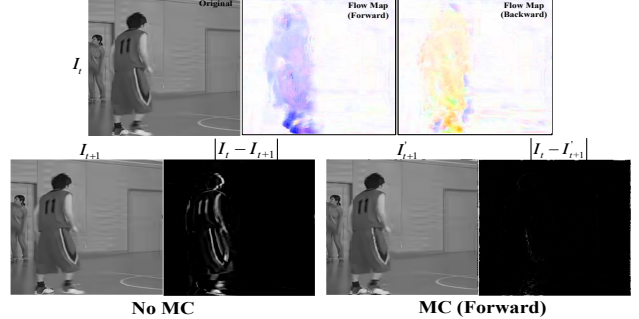
The proposed network consists of a temporal fusion subnet and a spatial detail enhancement subnet. Particularly, the temporal fusion subnet has two key modules, i.e., the motion compensation (MC) module and the fusion module. In this section, we focus on the design of these three modules.

#### 3.1. MC Module

The multi-frame video processing networks are normally built upon the fact that different observations of the same object or scene are probably available across frames of a video. As a result, content or scene, which are lost due to certain processing on the target frame, may be found in adjacent frames. Therefore, an intuitive idea is to enhance the compression quality of target frame by directly inputting multiple frames to the network. However, due to inter-frame motion, the interference information may be introduced to the network, especially for those scenes with drastic motion. To tackle this problem, we firstly employ a subnet to estimate and compensate the temporal motion across frames. Then, the compensated adjacent frames are used to enhance the quality of target frame.

In [7], Caballero *et al.* proposed the spatial transformer motion compensation (STMC) for video super-resolution. The basic idea of STMC is to predict the optical flow of adjacent frames to current frame by multi-scale down-sampling network. Suppose  $I_t$  and  $I_{t+1}$  are two consecutive frames, the optical flow related to adjacent frame  $I_{t+1}$ , whose reference frame is  $I_t$ , is a function of motion parameter  $\theta_{\Delta,t+1}$ . This optical flow can be represented by two feature maps corresponding to displacements of the  $x$  and  $y$  dimensions, i.e.,  $\Delta_{t+1}^x$  and  $\Delta_{t+1}^y$ , as  $\Delta_{t+1} = (\Delta_{t+1}^x, \Delta_{t+1}^y; \theta_{\Delta,t+1})$ . Then, the compensated frame  $I'_{t+1}$  can be expressed as

$$I'_{t+1}(x, y) = \mathcal{I} \{ I_{t+1}(x + \Delta_{t+1}^x, y + \Delta_{t+1}^y) \}, \quad (1)$$



**Fig. 2.** Top: flow map estimated relating the original frame. Bottom: the consecutive frames without and with MC (No MC and MC).

where  $\mathcal{I}$  denotes the bilinear interpolation. STMC consists of a coarse ( $\times 4$ ) and a fine ( $\times 2$ ) scale optical flow estimation.

We make several modifications on STMC to adapt it to our proposed SDTS method. First, we employ the coarse-to-fine ( $\times 4$  and  $\times 2$ ) flow estimation modules to handle large scale motion. Also, we develop a flow estimation module without down-sampling processing to handle still scenes in the video. Therefore, the final compensated frame  $I'_{t+1}$  is obtained by warping the target frame with the total flow

$$I'_{t+1} = \mathcal{I} \left\{ I_{t+1} \left( \Delta_{t+1}^c, \Delta_{t+1}^f, \Delta_{t+1}^s \right) \right\}, \quad (2)$$

where  $\Delta_{t+1}^c$ ,  $\Delta_{t+1}^f$  and  $\Delta_{t+1}^s$  denote the coarse flow, fine flow and still scenes flow, respectively. Second, we find that motion compensation relies to a large extent on the accuracy of motion estimation. Therefore, the proposed MC module is firstly trained under the supervision of raw frames to get a more accurate motion estimation, then the whole network is jointly fine-tuned based on this MC module.

To verify the effectiveness of our proposed MC module, we present the error maps between two consecutive frames  $I_t$  and  $I_{t+1}$  in Fig. 2. One can see from Fig. 2 that using the proposed MC operation induces less error in the compensated frame, and our proposed MC method can well eliminate interference in the adjacent frame.

#### 3.2. Multi-frame Fusion Module

The CNN-based temporal information fusion methods have been proposed for various applications, which are mainly classified into early fusion [19], slow fusion [20] and 3D convolutions [21]. Early fusion is one of the most straightforward fusion methods, which collapses all temporal information in the first layer. Slow fusion partially merges temporal information in a hierarchical structure and it is slowly fused as information progresses through the network, this fusion method has shown better performance than early fusion for some video applications [7, 20]. Therefore, we adopt the slow fusion mode as the temporal information fusion method in the SDTS and more details can be found in Fig. 3.

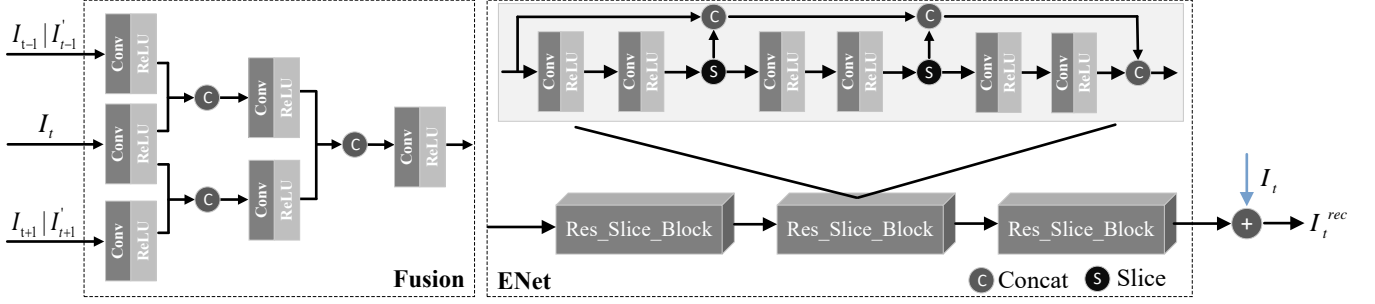


Fig. 3. The temporal fusion module and spatial detail enhancement subnet.

### 3.3. Enhancement subnet (ENet)

The enhancement subnet (ENet) is used to reduce the compression artifacts and enhance the reconstruction effect of target frame in our work. The experimental results in [22] and [23] demonstrate that adaptively recalibrating the responses of channel-wise features with coarse-to-fine structure can improve the representation of the network. Therefore, we construct our ENet with a series of coarse-to-fine residual slice blocks (Res.Slice.Block), as shown in Fig. 3. Specifically, only a part of the previous features are delivered to the following modules in each Res.Slice.Block to extract useful information progressively. The local short-path information and the local long-path information are aggregated by concatenation. The slice and concatenation in the Res.Slice.Block are used to control how much the useful information in current state will be reserved and delivered to the next unit. When the weights of both the operators are close to zeros, the information delivered from the previous state will be ignored by the current state. Conversely, more useful information in previous state will be delivered to current state.

### 3.4. Training Strategy

**Phase 1** The MC module is firstly trained under the supervision of raw frames  $I_0^R$  to get more accurate optical flow information. The loss of MC module can be written as

$$\mathcal{L}_{ME} = \sum_{i=-T}^T \|\mathcal{I}(I_{0 \rightarrow i}^R; \Delta_i^R) - I_0^R\|^2. \quad (3)$$

**Phase 2** We use Euclidean loss between the reconstructed target frame  $I_0^{Rec}$  and the ground truth  $I_0^H$  to train the quality enhancement subnet,

$$\mathcal{L}_{ENet} = \sum_{i=-T}^T \|I_0^H - I_0^{Rec(i)}\|^2. \quad (4)$$

**Phase 3** We finally fine-tuned the SDTS network by a joint loss function,

$$\mathcal{L} = \mathcal{L}_{ME} + \lambda_2 \mathcal{L}_{ENet}, \quad (5)$$

where  $\lambda_2$  is the weighting factor that balances the loss terms.

## 4. EXPERIMENT

We implement our SDTS framework on TensorFlow platform [24]. All the experiments are conducted on a PC with Intel Xeon E5 CPU and NVidia GeForce GTX 1080Ti GPU. We conduct all experiments on the same dataset to make a fair comparison between various methods, and all compared methods are retrained over the training dataset according to authors' recommended parameters.

**Data Preparation** The training and test sequences in the experiment are compressed in the common test conditions (CTCs) [26] by VVC reference software, VTM3.0, under Low-Delay P (LD) configuration. We specify the Quantization Parameters (QPs) to 32 and 37, respectively. When training the SDTS models, in each video clip, we randomly select the raw frame, its corresponding decoded target frame and the adjacent frames to form the training pairs.

**Model Training** All the proposed models are trained following the same protocol and share similar hyper-parameters. Filter sizes are set to  $3 \times 3$ , and all non-linearities are rectified linear units except for the output layer, which uses a linear activation. During training, we use a mini-batch size of 8. To minimize the loss functions of (5),  $\lambda_2$  is empirically set to 0.01, we employ Adam optimizer [27] with a start learning rate of  $1e-4$ , decay the learning rate with a power of 10 at the  $10^{th}$  epochs, and terminate training at 30 epochs. To save training time, we first train the model at QP 37 from scratch, and the model at QP 32 is fine-tuned from it.

In VVC, the distance between two HQFs that encoded under the LD configuration is normally less than five frames, such a short distance indicates that there exist high correlations among adjacent frames. As mentioned earlier, the LQFs may be enhanced using the adjacent HQFs. In addition, since the quality fluctuations across compressed video frames under LD configuration, it is difficult to enhance all video frames by utilizing a single model. In this work, we train a separate model for LQFs and HQFs, respectively, to enhance the quality of VVC compressed video. Both the trained models for the quality enhancement of LQFs and HQFs are proposed by taking advantage of the nearest adjacent HQFs of the corresponding target frame.

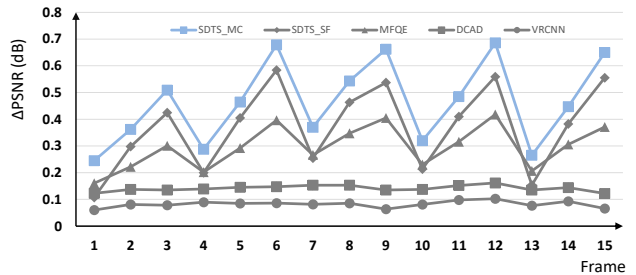
**Table 1.** Comparisons of different methods on  $\Delta$ PSNR (dB) over VTM3.0 baseline at QPs 37 and 32

QP	Class	Seq.	VRCNN [10]	DCAD [13]	MFQE [6]	SDTS (SF)	SDTS (MC)
37	B	Kimono1	0.0656	0.0743	0.1609	0.1793	0.2687
		ParkScene	0.1013	0.1285	0.2404	0.2485	0.3664
		Cactus	0.0832	0.0746	0.1548	0.2232	0.2799
		BasketballDrive	0.0452	0.0711	-0.0421	-0.0571	-0.0376
		BQTerrace	0.1249	0.1932	0.1953	0.2375	0.2882
	C	RaceHorsesC	0.0671	0.0691	0.0583	0.0738	0.1412
		BQMall	0.1029	0.1232	0.2155	0.3396	0.3835
		PartyScene	0.0634	0.0582	0.1206	0.2600	0.3044
		BasketballDrill	0.0704	0.1006	0.1056	0.1705	0.1984
	D	RaceHorses	0.1407	0.1763	0.2192	0.2071	0.3567
		BQSquare	0.1893	0.1947	0.1895	0.4081	0.4460
		BlowingBubbles	0.0833	0.1206	0.2371	0.3031	0.3754
		BasketballPass	0.1251	0.1455	0.3027	0.3866	0.4882
	E	FourPeople	0.2185	0.2232	0.2997	0.4593	0.4761
		Johnny	0.2029	0.2105	0.2906	0.3957	0.4256
		KristenAndSara	0.1943	0.2013	0.2934	0.4272	0.4577
		<b>Average</b>	<b>0.1174</b>	<b>0.1353</b>	<b>0.1939</b>	<b>0.2664</b>	<b>0.3262</b>
32	<b>Average</b>	<b>0.1004</b>	<b>0.1251</b>	<b>0.1974</b>	<b>0.2703</b>	<b>0.2871</b>	

#### 4.1. Quantitative Evaluation

To verify the performance of the proposed SDTS-based method, we evaluate the performance of our SDTS method in terms of  $\Delta$ PSNR, which measures the PSNR difference between the enhanced and the original compressed frame. We compare our SDTS-based method with some state-of-the-art algorithms, that is, VRCNN [10], DCAD [13] and MFQE [6]. Particularly, VRCNN and DCAD are single-frame based methods, while MFQE is a multi-frame based video quality enhancement method. In addition, in order to verify the effect of MC module in the temporal fusion subnet, we also retrained a model that disables the MC module as a comparison. The temporal fusion component of the comparison model is only implemented by slow fusion, and the corresponding retrained model is denoted as SDTS.SF. The trained model with MC module is denoted SDTS.MC.

Table 1 presents the  $\Delta$ PSNR results of each test sequence at QPs 37 and 32. It can be seen from Table 1 that our SDTS method outperforms (on average) the other methods. Specifically, the highest  $\Delta$ PSNR of our SDTS reaches 0.4882dB for MC mode at QP 37, the averaged  $\Delta$ PSNR gains of our SDTS method is 0.3262dB and 0.2664dB for MC and SF modes, respectively, which are much higher than that of MFQE method (0.1939dB), the state-of-the-art method. Besides, when compared with VRCNN and DCAD methods, the SDTS method can even achieve a much higher PSNR gain. As shown in Ta-



**Fig. 4.** Quality fluctuation of “BlowingBubbles” at QP 37.

ble 1, this result has a similar trend at QP 32. Based on these results, we can conclude that the spatial prior information and multi-frame temporal information play important roles in the quality enhancement for VVC compressed videos.

#### 4.2. Quality Fluctuation

We also compare the quality fluctuation of compressed video between various methods. As shown in Fig. 4, we provide the quantified  $\Delta$ PSNR results for 15 consecutive frames of the test video “BlowingBubbles”. Although we only present the results of this test sequence in Fig. 4, we find that all other test sequences have similar results. From Fig. 4, one can see that the  $\Delta$ PSNR curve of our SDTS method is always over the  $\Delta$ PSNR curves of comparison methods, which indicates that our method can reach a higher  $\Delta$ PSNR gain for each single frame than the comparison methods. Therefore, our proposed SDTS method is effective to mitigate the quality fluctuation of VVC compressed video, as well as enhancing the quality of compressed video.

### 5. CONCLUSIONS

We propose a novel CNN-based method to enhance the VVC compressed videos by jointly exploiting spatial details and temporal structure. Our proposed method, i.e. the SDTS-based network, consists of a temporal information fusion subnet and a spatial detail enhancement subnet. The former subnet is utilized to estimate and compensate the temporal motion across frames, and the latter one is employed to enhance the reconstruction quality of the VVC compressed video. Experimental results demonstrate that the proposed SDTS-based method achieves the state-of-the-art performance.

**Acknowledgement** This work was partially supported by the Applied Basic Research Program of Sichuan Province under Grant 2019YJ0163, by the National Natural Science Foundation of China under Grant 61672134, 61701310, by the Free Exploration Grant for Basic Research of Shenzhen City JCYJ20180305124209486, and by the Fundamental Research Funds for Central Universities of China under Grant ZYGX2016J038.

## 6. REFERENCES

- [1] J.-R. Ohm and G. J. Sullivan, "Versatile video coding-towards the next generation of video compression," in *PCS*, 2018.
- [2] G. J. Sullivan, J. R. Ohm, W.-J. Han, and T. Wiegand, "Overview of the high efficiency video coding (HEVC) standard," *IEEE Trans. on Circuits and Systems for Video Technology*, vol. 22, no. 12, pp. 1649–1668, 2012.
- [3] F. Brandi, R. de Queiroz, and D. Mukherjee, "Super resolution of video using key frames," in *ISCAS*, 2008, pp. 1608–1611.
- [4] E. M. Hung, R. L. de Queiroz, F. Brandi, K. F. de Oliveira, and D. Mukherjee, "Video super-resolution using codebooks derived from key-frames," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 22, no. 9, pp. 1321–1331, 2012.
- [5] B. C. Song, S.-C. Jeong, and Y. Choi, "Video super-resolution algorithm using bi-directional overlapped block motion compensation and on-the-fly dictionary training," *IEEE Trans. on Circuits and Systems for Video Technology*, vol. 21, pp. 274–285, 2011.
- [6] R. Yang, M. Xu, Z. Wang, and T. Li, "Multi-frame quality enhancement for compressed video," in *CVPR*, 2018, pp. 6664 – 6673.
- [7] J. Caballero, C. Ledig, A. Aitken, A. Acosta, J. Totz, Z. Wang, and W. Shi, "Real-time video super-resolution with spatio-temporal networks and motion compensation," in *CVPR*, 2017.
- [8] X. Tao, H. Gao, Y. Wang, X. Shen, J. Wang, and J. Jia, "Scale-recurrent network for deep image deblurring," in *CVPR*, 2018.
- [9] C. Wang, H. Huang, X. Han, and J. Wang, "Video inpainting by jointly learning temporal structure and spatial details," *preprint arXiv:1806.08482*, 2018.
- [10] Y. Dai, D. Liu, and F. Wu, "A convolutional neural network approach for post-processing in HEVC intra coding," in *MMM*, 2017, pp. 28 – 39.
- [11] C. Dong, Y. Deng, C. C. Loy, and X. Tang, "Compression artifacts reduction by a deep convolutional network," in *ICCV*, 2015, pp. 576 – 584.
- [12] Y. Tai, J. Yang, X. Liu, and C. Xu, "MemNet: A persistent memory network for image restoration," in *ICCV*, 2017.
- [13] T. Wang, M. Chen, and H. Chao, "A novel deep learning-based method of improving coding efficiency from the decoder-end for HEVC," in *DCC*, 2017.
- [14] Z. Wang, D. Liu, S. Chang, Q. Ling, Y. Yang, and T. S. Huang, "D3: Deep dual-domain based fast restoration of jpeg-compressed images," in *CVPR*, 2016, pp. 2764 – 2772.
- [15] X. Meng, X. Deng, S. Zhu, S. Liu, C. Wang, C. Chen, and B. Zeng, "MGANet: A robust model for quality enhancement of compressed video," *arXiv:1811.09150*, pp. 1–12, 2018.
- [16] 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.
- [17] X. He, Q. Hu, X. Zhang, C. Zhang, W. Lin, and X. Han, "Enhancing HEVC compressed videos with a partition-masked convolutional neural network," in *ICIP*, 2018, pp. 216–220.
- [18] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016, pp. 770 – 778.
- [19] A. Kappeler, S. Yoo, Q. Dai, and A. K. Katsaggelos, "Video super-resolution with convolutional neural networks," *IEEE Transactions on Computational Imaging*, vol. 2, no. 2, pp. 109 – 122, 2016.
- [20] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei, "Large-scale video classification with convolutional neural networks," in *CVPR*, 2014.
- [21] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri, "Learning spatiotemporal features with 3D convolutional networks," in *CVPR*, 2015, pp. 4489–4497.
- [22] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *CVPR*, 2018, pp. 7132 – 7141.
- [23] Z. Hui, X. Wang, and X. Gao, "Fast and accurate single image super-resolution via information distillation network," in *CVPR*, 2018, pp. 723–731.
- [24] M. Abadi, A. Agarwal, and Paul Barham et al., "TensorFlow: Large-scale machine learning on heterogeneous systems," 2015, Software available from tensorflow.org.
- [25] L. Song, X. Tang, W. Zhang, X. Yang, and P. Xia, "The sju 4k video sequence dataset," in *QoMEX*, 2013.
- [26] J. Boyce, K. Suehring, X. Li, and V. Seregin, "JVET common test conditions and software reference configurations," *JVET-J1010, ITU-T SG16*, 2018.
- [27] D. Kingma and B. Jimmy, "Adam: A method for stochastic optimization," in *ICLR*, 2014.