

CPGA: Coding Priors-Guided Aggregation Network for Compressed Video Quality Enhancement

Qiang Zhu¹, Jinhua Hao², Yukang Ding², Yu Liu¹, Qiao Mo¹, Ming Sun², Chao Zhou², Shuyuan Zhu^{1,*}

¹University of Electronic Science and Technology of China

²Kuaishou Technology

{zhuqiang@std., 202211012315@std., eezsy@}uestc.edu.cn, mqiao568@gmail.com

{haojinhua, dingyukang, sunming03, zhouchao}@kuaishou.com

Abstract

Recently, numerous approaches have achieved notable success in compressed video quality enhancement (VQE). However, these methods usually ignore the utilization of valuable coding priors inherently embedded in compressed videos, such as motion vectors and residual frames, which carry abundant temporal and spatial information. To remedy this problem, we propose the **Coding Priors-Guided Aggregation (CPGA)** network to utilize temporal and spatial information from coding priors. The CPGA mainly consists of an inter-frame temporal aggregation (ITA) module and a multi-scale non-local aggregation (MNA) module. Specifically, the ITA module aggregates temporal information from consecutive frames and coding priors, while the MNA module globally captures spatial information guided by residual frames. In addition, to facilitate research in VQE task, we newly construct the **Video Coding Priors (VCP)** dataset, comprising 300 videos with various coding priors extracted from corresponding bitstreams. It remedies the shortage of previous datasets on the lack of coding information. Experimental results demonstrate the superiority of our method compared to existing state-of-the-art methods. The code and dataset will be released at <https://github.com/CPGA/CPGA.git>.

1. Introduction

Video contains and communicates perceptual information derived from the real world. With the growth of the Internet, it has witnessed explosive growth in video content in digital network traffic [33]. When transmitting videos over the Internet with limited bandwidth, efficient video codecs, such as H.264/AVC [32], H.265/HEVC [27], AV1 [4] and H.266/VVC [2, 11] are widely used to save the coding bitrate. However, due to the quantization in the lossy video

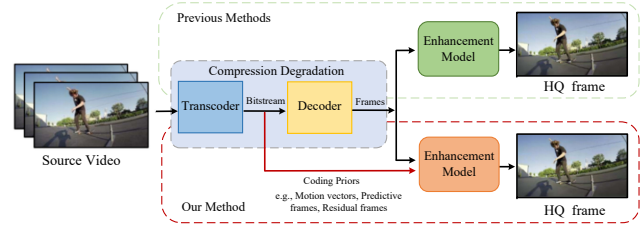


Figure 1. Comparison between previous methods and our CPGA. Compared with previous methods, the coding priors are extracted from the bitstream to input into our enhancement model for VQE.

coding, compressed videos inevitably exhibit compression artifacts, significantly diminishing the quality of experience [1, 15, 26, 28, 42, 48, 49]. Moreover, compression artifacts may degrade the accuracy of vision tasks applied to videos, including object detection [20, 53] and action recognition [9, 52]. Hence, improving the quality of compressed videos is of importance and practical significance.

In the past decade, extensive methods have emerged for compressed video quality enhancement (VQE), especially after the rise of deep learning. Inspired by image restoration [23, 24, 54] task for improving the quality of degraded images, the single-frame-based methods [5, 7, 21, 30] are firstly proposed to achieve substantial breakthroughs by learning the inverse mapping from spatial information of low-quality (LQ) and high-quality (HQ) frames. Moreover, multi-frame-based approaches [6, 10, 16, 17, 34, 39, 40, 51] devote to exploiting the temporal correlation between multiple consecutive LQ frames, consequently yielding superior performance on details restoration over single-frame-based methods in VQE task. To leverage the temporal information, researchers commonly utilize optical flow [10, 40] or employ deformable convolutions [6, 17, 51] to align adjacent frames toward the HQ frame to be reconstructed. However, the introduction of optical flow estimation inevitably increases the computational overhead, and inaccurate flow estimation and alignment may degrade the final quality [6].

*Corresponding author

It is noteworthy that compressed videos inherently encompass coding information, such as motion vector (MV), predictive frame, residual frame, and partition map, collectively known as coding priors. MVs represent the temporal information between adjacent frames, predictive frames provide aligned results from the decoding process and the residual frame represents the difference between the current frame and the corresponding predictive frame. These coding priors can be extracted from the bitstream, characterizing explicit temporal and spatial information of compressed videos.

Recently, various high-level vision tasks [43, 50] and low-level vision tasks [14, 29] have benefited from such coding priors. Meanwhile, a few approaches leveraging coding priors have yielded substantial success in compressed video super-resolution (VSR) [3, 31, 44]. Notably, compressed domain deep video super-resolution (CD-VSR) [3] establishes a pioneering dataset with coding priors for the VSR task. Despite the commendable performance achieved in these tasks, the full and effective utilization of valuable coding priors within compressed videos remains inadequately addressed in the VQE task. Additionally, existing VQE datasets typically comprise solely raw and compressed videos, lacking explicit provision of coding priors. Recognizing this, we are inspired to explore the potential of coding priors for further improving performance in VQE task.

In this paper, we firstly construct a new Video Coding Priors (VCP) dataset, comprising 300 raw videos compressed by various HEVC configurations and corresponding coding priors extracted from the bitstream, including MVs, predictive frames and residual frames. VCP dataset can remedy the shortage of previous datasets on the lack of coding priors. Based on VCP dataset, we propose a novel coding priors-guided aggregation network, named CPGA. The CPGA consists of three modules: the inter-frame temporal aggregation (ITA) module, the multi-scale non-local aggregation (MNA) module and the quality enhancement (QE) module. Specifically, the ITA module explores the inter-frame correlations among the multiple compressed frames with the guidance of MVs and predictive frames to generate effective temporal features. Then, the MNA module is designed to globally capture the spatial correlations within the feature. It obtains the spatially-aggregated features with the guidance of the current residual frames. After that, the QE module is constructed to enhance the spatially-aggregated feature to generate the HQ video frame. We illustrate the comparison between previous methods and our method in Fig. 1. With the benefits of coding priors and effective feature aggregation modules, our CPGA outperforms previous methods on public testing sequences. The main contributions can be summarized as follows:

- We establish a compressed videos with coding priors

dataset for VQE task, named as VCP, which includes LQ sequences, HQ sequences and three coding priors extracted from bitstream (MV, predictive frames and residual frames). The dataset remedies the shortage of previous datasets without coding priors in VQE task.

- We propose a coding priors-guided aggregation (CPGA) network for VQE task. The CPGA composed of feature aggregation modules can efficiently achieve better temporal and spatial features representation with leveraging valuable coding priors in our dataset.
- Experimental results demonstrate that our method achieves a performance gain of more than 0.03dB compared to previous state-of-the-art methods and outperforms [17] by 10% on inference speed.

2. Related Work

2.1. Compressed Video Quality Enhancement

In the early years, several single-frame-based methods [5, 7, 12, 14, 21, 30, 38, 47] were introduced to leverage a LQ frame of compressed video to generate a HQ frame, aiming at enhancing the quality of compressed videos. For instance, the artifact reduction convolutional neural network [7] employed a 4-layer CNN to reduce compression artifacts for Joint Photographic Experts Group (JPEG) images. Subsequently, the residual non-local attention network [47] was introduced to capture long-range dependencies between pixels for JPEG image enhancement. To leverage prior information from compressed video, a partition-aware convolutional neural network [14] was employed to utilize coding unit information from partition images to enhance the quality of compressed frames.

Recently, several multi-frame-based methods have been proposed for enhancing the quality of compressed videos by utilizing information from multiple frames. In these methods, optical flow is employed as motion information to align adjacent frames for quality enhancement [10, 40]. For instance, the multi-frame quality enhancement (MFQE) network [40] initially used optical flow to align adjacent compressed frames, achieving quality enhancement for compressed video. Subsequently, MFQE2.0 [10] introduced a multi-scale approach, batch normalization and dense connection to further enhance performance. However, in scenarios with significant motion, even a small alignment error may lead to serious artifacts in aligned frames, resulting in poor quality of the composed HR frame. In addition to flow-based methods, deformable convolution-based methods have been proposed [6, 10, 16] to learn offsets from compressed frames to obtain aligned features for VQE. For instance, spatio-temporal deformable fusion (STDF) network [6] jointly predicts all deformable offsets for multiple frames based on deformable convolution network (DCN) [55] to achieve enhancement. Furthermore,

Table 1. Comparison between existing video enhancement datasets and our dataset.

Dataset	Type	Number	Resolution	Compression Settings	Coding Priors
LDV [37]	Training+Validation+Test	240	960 × 536	LDP at QP=37	-
LDV2.0 [41]	Training+Validation+Test	335	4K, 960 × 536	LDP at QP=37	-
MFQE2.0 [10]	Training	106	SIF, CIF, 720×480, 4CIF, 360p, 1080p and 2K	LDP at QPs=22,27,32,37,42	-
VCP (Ours)	Training	300	272×480, SIF, CIF, 640×480, 720×480, 4CIF, 360p, 720p, 1080p, 2K and 4K	LDB & RA at QPs=22,27,32,37	✓

recursive fusion and deformable spatiotemporal attention (RFDA) network [51] employs a recursive fusion scheme to exploit relevant information from multiple frames based on DCN [55] in a large temporal range. The spatio-temporal detail information retrieval (STDR) network [17] introduces a multi-path deformable alignment module to generate more accurate offsets by integrating the alignment features of different receptive fields, enabling better recovery of temporal detail information for generating HQ videos.

Although previous multi-frame-based methods achieve state-of-the-art performance, they ignore the coding information of compressed videos, limiting enhancement performance improvement. Such coding priors contain additional temporal and spatial information, and naturally embedded within the bitstream of compressed video. Recently, a few pioneering studies [3, 31, 44] have been proposed to incorporate coding priors in the VSR task. For instance, CD-VSR [3] designed a new framework to utilize the deep priors and coding priors to achieve the improvement of performance for VSR. Besides, the codec information assisted framework [44] utilized the motion vector to model the temporal relationships between adjacent frames and skip the redundant pixels based on the residual frame to achieve high-efficient VSR. Inspired by these works, we have constructed a compressed video quality enhancement dataset that includes LQ frames and their coding priors. Based on this dataset, we design a coding priors-guided aggregation network for VQE.

2.2. Video Enhancement Datasets

In the past decade, numerous datasets [10, 18, 22, 35, 37, 41] have been developed for compressed video quality enhancement and they consist of HQ sequences and corresponding compression configurations, such as Low Delay P (LDP) configuration and different quantization parameters (QPs). We summarize the characteristics of some representative datasets in Tab. 1. Although compressed video quality enhancement methods developed based on these datasets [10, 37, 41] have achieved notable success, their applicability has been largely limited in exploring video representations to further improve enhancement performance. Moreover, some coding priors (e.g., MVs, predictive frames, and

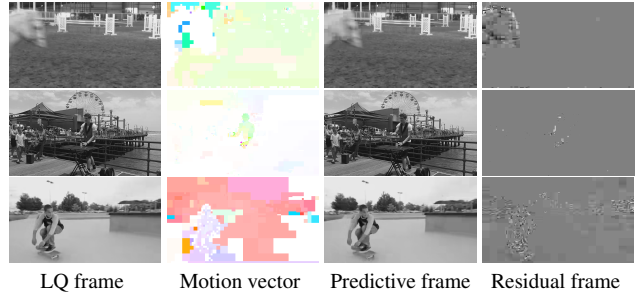


Figure 2. Some examples of LQ frames, motion vectors, predictive frames and residual frames in proposed VCP dataset.

residual frames) contain rich temporal and spatial information extracted from compressed videos, which can assist in constructing HQ videos. To efficiently explore the potential performance advantages offered by the coding priors of compressed videos, we propose our VCP dataset to accommodate interactions between compressed video quality enhancement and video coding.

3. Our dataset

Data collection. In the compressed video quality enhancement task, existing datasets [18, 37, 40] typically include only raw and compressed videos, lacking the provision of coding priors, which limits design of our network. Herein, we collect 300 videos from CDVL [13] and original VIMEO videos [36]. The raw videos are selected across diverse types of content, such as wilds, urban, daily routines and professional sports. Accordingly, 300 video clips are extracted from these videos and each of them consists of 48 frames. All LQ sequences are compressed with HEVC software HM 16.25 [27]. In particular, we adopt two configurations, i.e., Low Delay B (LDB) and Random Access (RA), and employ four QPs: 22, 27, 32 and 37, greatly facilitating the study of joint video coding and enhancement. Additionally, we extract three types of coding priors embedded in the bitstream of compressed videos, which include MVs, predictive frames and residual frames.

Dataset analysis. Some examples from our VCP dataset, as shown in Fig. 2, cover a wide range of real-world scenarios. This diversity allows for a comprehensive evaluation of enhancement performance across different appli-

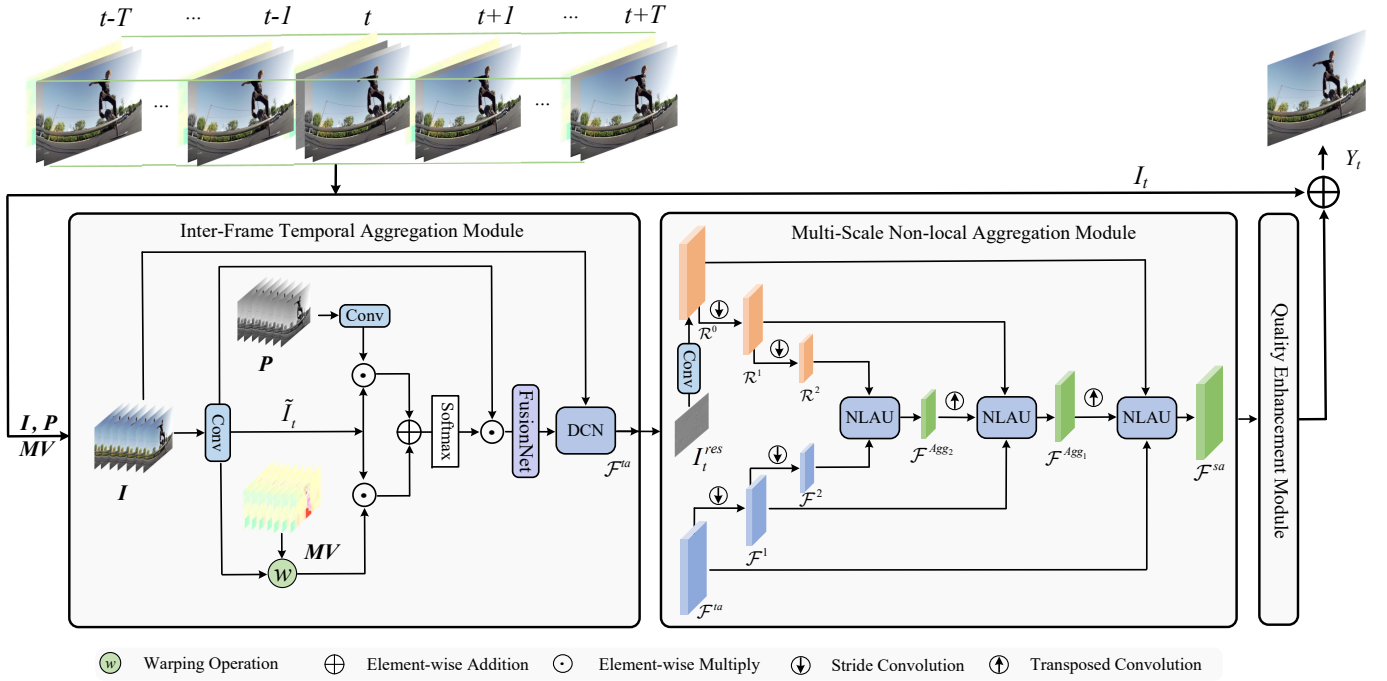


Figure 3. The architecture of CPGA. The MVs and predictive frames are fed into the ITA module to obtain the temporally-aggregated feature, and current residual frame is utilized in the MNA module to obtain the spatially-aggregated feature.

cations, including autonomous driving, video surveillance and photo editing. In addition to the diversity of scenarios, our VCP dataset includes various videos with different resolutions from 272×480 to 4K. Two coding configurations and four QPs are adopted to generate the different degradation compressed videos from raw videos, which facilitates the evaluation of the effectiveness of compressed video quality enhancement methods. Furthermore, the corresponding three coding priors of compressed videos also were extracted, i.e., MVs, predictive frames and residual frames, which provides the facilitation for the application of coding priors in VQE task. We summarize the characteristics of our dataset in Tab. 1. This comprehensive dataset facilitates the evaluation of the effectiveness of compressed video quality enhancement methods.

4. Methodology

4.1. Overall Framework

The architecture of our CPGA network is illustrated in Fig. 3. It comprises three key modules: the inter-frame temporal aggregation (ITA) module, the multi-scale non-local aggregation (MNA) module and the quality enhancement (QE) module. Firstly, ITA module conducts spatial feature extraction and temporal feature fusion on consecutive LQ frames and corresponding two coding priors, i.e., MVs and predictive frames. ITA module exploits temporal correlations of adjacent frames to generate temporally-aggregated features under the guidance of these two coding priors. Subsequently, the temporally-aggregated features and the current residual frame are input into the MNA

module, where non-local multi-scale features are integrated, resulting in the generation of spatially-aggregated features. Finally, the generated features are fed into the QE module to derive the ultimately enhanced feature. The details of the proposed ITA, MNA and QE modules are explained in the following Secs. 4.2 to 4.4.

4.2. Inter-frame Temporal Aggregation Module

ITA module aims to effectively align features of adjacent frames, it incorporates two coding priors, i.e., MVs and predictive frames, to guide the fusion of features of LQ frames, thus generating the temporally-aggregated features.

Given a sequence of $2T + 1$ consecutive LQ frames $I = I_{[t-T:t+T]}$, where I_t represents the target frame to be enhanced and the others are reference frames. Corresponding predictive frames and motion vectors are denoted as $P = P_{[t-T:t+T]}$ and $MV = MV_{[t-T+1:t+T]}$, note that there are $2T$ MVs in $2T+1$ frames. We utilize a convolution layer with 3×3 kernel size on $I_{[t-T:t+T]}$ and $P_{[t-T:t+T]}$ to extract their features. The extracted features are denoted as $\tilde{I}_{[t-T:t+T]}$ and $\tilde{P}_{[t-T:t+T]}$. In addition, leveraging the motion vectors $MV_{[t-T+1:t+T]}$, the features of the LQ frames are aligned via a *Warp* operation [8] and obtain the aligned features $\tilde{F}_{[t-T+1:t+T]}$. Since the first frame does not have the corresponding MV from the forward frame, we reuse \tilde{I}_{t-T} to represent the first aligned feature \tilde{F}_{t-T} . Usually, the predictive frames are aligned results from the video decoding process while aligned features are the post-process results from motion vectors. Thus, we use these two aligned results to explore the inter-frame correlations among frames for generating temporally-aggregated features.

Specifically, the inter-frame correlations between current LQ frame feature \tilde{I}_t , features of predictive frames $\tilde{P}_{[t-T:t+T]}$ and aligned features $\tilde{F}_{[t-T:t+T]}$ are explored via

$$\hat{P}_{[t-T:t+T]} = \tilde{I}_t \odot \tilde{P}_{[t-T:t+T]} \quad (1)$$

and

$$\hat{F}_{[t-T:t+T]} = \tilde{I}_t \odot \tilde{F}_{[t-T:t+T]}, \quad (2)$$

where \odot is an element-wise multiply operation. Then, these two temporal features are added and multiplied with the features of LQ frames for compensating the temporal information of LQ frames to obtain the temporally-compensated features, as shown in Eq. (3)

$$\mathcal{F}_i^c = \tilde{I}_i \odot \sigma(\hat{P}_i + \hat{F}_i), \quad (3)$$

here $i \in [t-T, t+T]$ and σ is softmax function. After that, we use a *FusionNet* consisting of a convolution layer with 3×3 kernel size and a similar Unet[25] from STDF [6], to obtain the fused feature via fusing temporally-compensated features, as shown in Eq. (4)

$$\mathcal{F}_t^f = \text{FusionNet}(\mathcal{F}_{[t-T:t+T]}^c). \quad (4)$$

To further aggregate temporal information, we employ a DCN [55] to aggregate the temporal feature \mathcal{F}_t^f to finally obtain the temporally-aggregated feature \mathcal{F}^{ta} .

4.3. Multi-Scale Non-Local Aggregation Module

In video coding, residual frames refer to the difference between the frame predicted from the previous frame and the current frame. The difference is typically concentrated in essential regions such as edges and contour areas of frames. To fully aggregate the spatial information, we design the MNA module to globally explore the spatial information of features with the guidance of the current residual frame at three scales. The structure of the MNA module is shown in Fig. 3.

Specifically, the current residual frame is firstly fed into a convolution layer with 3×3 kernel size to obtain its feature, denoted as \mathcal{R}^0 . A non-local aggregation unit (NLAU) is designed to globally explore spatial information of features with the guidance of the feature \mathcal{R}^0 . Then, NLAU is adopted at three different scales ($s=0, 1, 2$) to progressively aggregate features to generate spatial-aggregated feature. We denote temporally-aggregated feature \mathcal{F}^{ta} as feature \mathcal{F}^0 at scale $s=0$ in the MNA module for spatial aggregation. In detail, feature \mathcal{F}^0 and \mathcal{R}^0 are sent into the convolution layer with a stride size of 2 to obtain corresponding downscaled features $\mathcal{F}^1, \mathcal{F}^2$ and $\mathcal{R}^1, \mathcal{R}^2$.

At scale $s=2$, feature $\mathcal{F}^2, \mathcal{R}^2$ are firstly fed into the NLAU without the previous aggregated feature to explore the spatial correlations within frames for aggregating them to obtain feature \mathcal{F}^{Agg2} . After that, feature \mathcal{F}^{Agg2} is upsampled by a transposed convolution layer with a stride size

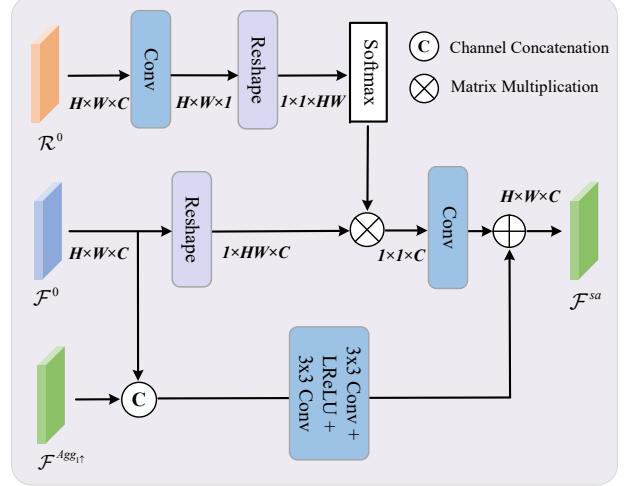


Figure 4. The structure of non-local aggregation unit (NLAU).

of 2 to obtain the upsampled feature $\mathcal{F}^{Agg2\uparrow}$. Subsequently, the feature $\mathcal{F}^{Agg2\uparrow}, \mathcal{F}^1, \mathcal{R}^1$ are fed into the NLAU to obtain the spatially-aggregated feature \mathcal{F}^{Agg1} at scale $s=1$. Finally, $\mathcal{F}^0, \mathcal{R}^0$ and the feature $\mathcal{F}^{Agg1\uparrow}$ are fed into NLAU to obtain the final spatially-aggregated feature \mathcal{F}^{sa} . This process is formulated as

$$\mathcal{F}^{Aggs} = \begin{cases} \text{NLAU}(\mathcal{F}^s, \mathcal{R}^s, \mathcal{F}^{Aggs+1\uparrow}), & s = 0, 1 \\ \text{NLAU}(\mathcal{F}^s, \mathcal{R}^s), & s = 2. \end{cases} \quad (5)$$

We illustrate the structure of NLAU at scale $s=0$ in Fig. 4. In NLAU, feature \mathcal{F}^0 is transformed with the guidance of feature \mathcal{R}^0 to achieve the spatial aggregation within feature to obtain the spatial feature. After that, the upsampled aggregated feature $\mathcal{F}^{Agg1\uparrow}$ and feature \mathcal{F}^0 are fused to maintain the previous aggregated information. Finally, the fused feature and the spatial feature are added to obtain the spatially-aggregated feature \mathcal{F}^{sa} .

4.4. Quality Enhancement Module

To further enhance spatially-aggregated feature \mathcal{F}^{sa} for the composition of the high-quality frame, we first employ two convolution layers with 3×3 kernel size followed by the LeakyReLU activate function and then develop G number of shift channel attention blocks (SCAB) to construct our QE module. The structure of the QE module and SCAB are illustrated in Fig. 5.

In SCAB, we embed the partial channel shifting operation [45] in the front of the channel attention block (CAB) [46] to enlarge the receptive field. Specifically, we choose the γ proportion feature in the middle position of the feature to use the partial channel shifting operations to enlarge the receptive field for feature enhancement without increasing parameters and complexity. We design two partial channel shifting operations with different directions, i.e., partial channel shifting operation along the horizontal direction

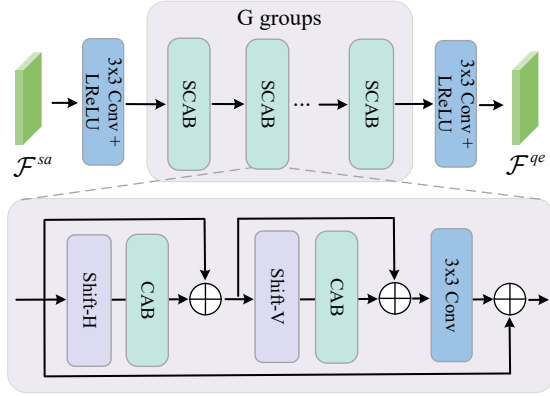


Figure 5. The structure of the quality enhancement (QE) module and the shift channel attention block (SCAB).

(*Shift-H*) with \hat{h} pixel and partial channel shifting operation along the vertical direction (*Shift-V*) with \hat{w} pixel. The *Shift-H* operation is first adopted to enlarge the receptive field along the horizontal direction at the front of CAB. Subsequently, the *Shift-V* operation is adopted to enlarge the receptive field along the vertical direction at the front of the CAB. Based on these two shift operations and CABs, the SCAB is designed to construct our QE module for feature enhancement.

5. Experiments

5.1. Experimental Settings

We adopt our VCP dataset as our training dataset. Following [6, 16, 17, 51], 18 standard testing sequences from common test conditions of JCT-VC [19] are used as our testing sequences. All testing sequences are processed under LDB and RA configurations with four QPs: 22, 27, 32, and 37, to generate the compressed videos and corresponding coding priors. The peak signal-to-noise ratio (PSNR) and structural similarity index (SSIM) [41] are adopted for performance evaluation. The enhanced results are evaluated in terms of Δ PSNR and Δ SSIM, which measures the PSNR/SSIM gap between the enhanced and the compressed videos. These two metrics are offered over the Y channel of YCbCr space as previous works [6, 16, 17, 51] did.

5.2. Implementation Details

In the training phase, we randomly crop 128×128 clips from raw videos, compressed videos and their corresponding coding priors as training samples with setting the batch size to 32. Each video clip contains 7 frames, i.e., $T=3$. We also employ data augmentation strategies (i.e., selection or flipping) to expand our dataset. We train all the models using Adam [55] optimizer with $\beta_1=0.9$, $\beta_2=0.999$, the compensation factor $\varepsilon = 1 \times 10^{-8}$ and the learning rate is initially set to 1×10^{-4} . We adopt the Charbonnier loss [30]

to train our proposed model. For setting of our model, we set $G=2$ in the QE module, $\hat{h} = \hat{w} = 2$ and γ is set to $1/8$ in SCAB. The kernel size of deformable convolution in the ITA module is set to 3×3 .

5.3. Comparison with State-of-the-art Methods

To demonstrate the effectiveness of our method, we compare our method with five multi-frame-based compressed video quality enhancement methods, i.e., MFQE [40], STDF-R3L [6], RFDA [51], coarse-to-fine spatio-temporal information fusion (CF-STIF) [16] and STDR [17]. We re-train these methods on our dataset using their settings, the results of these methods under LDB and RA configurations are provided in Tab. 2 and Tab. 3, respectively.

Quantitative Results. We present the quantitative results in terms of average Δ PSNR and Δ SSIM in Tab. 2 and Tab. 3. It can be observed from Tab. 2 and Tab. 3 that our method consistently outperforms all the compared methods, which demonstrates that our method achieves the state-of-the-art performance under LDB and RA configurations with all different QPs. Specifically, our method outperforms STDF-R3L and RFDA by 0.13dB and 0.08dB in terms of average Δ PSNR over 18 standard testing sequences under LDB configuration at QP = 37. Besides, compared with a state-of-the-art method, i.e., STDR, our method also achieves an improvement of 0.03dB under LDB configuration at QP = 37. Moreover, our model improves about 0.02dB-0.05dB in terms of average Δ PSNR under RA configuration.

Quality Fluctuation. Quality fluctuation is another observable measurement for the overall quality of enhanced videos. Drastic quality fluctuation of frames accounts for severe texture shaking and degradation of the quality of experience. We provide the PSNR curves of the HEVC, STDF-R3L, STDR and our method on a testing sequence, i.e., BasketballPass, in Fig. 6. It can be seen in Fig. 6 that our method has not only more improvement in performance but also smaller fluctuations.

Qualitative Results. We present the qualitative results of the compared methods and our method under different coding configurations and QPs in Fig. 7. It is found from Fig. 7 that our method obviously reduces more compression artifacts and generates better visual results.

5.3.1 Model Complexity

As described in [6, 51], the evaluation of model complexity involves considering both the model parameters and the processed frames per second (FPS). The corresponding results are given in Tab. 4. It is found from Tab. 4 that our method outperforms CF-STIF in terms of both parameters and FPS. The parameters of our CPGA model are similar to STDR, but CPGA achieves a higher inference speed. Overall, our CPGA demonstrates a superior balance between performance and efficiency.

Table 2. Quantitative comparison in terms of Δ PSNR (dB) and Δ SSIM ($\times 10^{-2}$) under LDB configuration. **Red** text indicates the best and **blue** text indicates the second best performance. * Video resolutions: Class A (2560 \times 1600), Class B (1920 \times 1080), Class C (832 \times 480), Class D (416 \times 240) and Class E (1280 \times 720).

QP	Class*	Sequence	MFQE		STDF-R3L		RFDA		CF-STIF		STDR		Ours	
			PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
37	A	Traffic	0.36	0.75	0.63	1.11	0.68	1.18	0.67	1.19	0.71	1.23	0.73	1.21
		PeopleOnStreet	0.56	1.13	1.05	1.89	1.15	1.98	1.12	2.00	1.17	2.04	1.27	2.12
		Kimono	0.21	0.55	0.64	1.22	0.68	1.37	0.68	1.28	0.72	1.32	0.79	1.43
	B	ParkScene	0.18	0.57	0.46	1.18	0.53	1.30	0.49	1.33	0.56	1.37	0.56	1.43
		Cactus	0.26	0.63	0.66	1.29	0.69	1.27	0.71	1.39	0.73	1.38	0.74	1.39
		BQTerrace	0.20	0.45	0.55	0.97	0.55	0.96	0.56	1.00	0.59	1.08	0.59	1.08
	C	BasketballDrive	0.27	0.56	0.65	1.09	0.70	1.15	0.70	1.21	0.74	1.28	0.78	1.26
		RaceHorses	0.22	0.57	0.49	1.29	0.51	1.26	0.56	1.32	0.58	1.36	0.62	1.57
		BQMall	0.21	0.77	0.84	1.76	0.90	1.79	0.90	1.88	0.97	1.96	0.97	1.91
	D	PartyScene	-0.12	0.34	0.53	1.76	0.58	1.77	0.62	1.84	0.63	1.91	0.61	2.02
		BasketballDrill	0.24	0.75	0.69	1.37	0.73	1.37	0.77	1.51	0.78	1.57	0.78	1.46
		RaceHorses	0.31	0.86	0.70	1.87	0.75	1.92	0.78	2.01	0.83	2.10	0.85	2.23
	E	BQSquare	-0.53	-0.29	0.69	1.14	0.82	1.36	0.83	1.41	0.87	1.41	0.89	1.43
		BlowingBubbles	0.12	0.77	0.57	2.07	0.68	2.28	0.66	2.31	0.68	2.41	0.67	2.43
		BasketballPass	0.27	0.79	0.82	1.76	0.86	1.75	0.92	1.99	0.97	2.08	0.98	2.10
Average	FourPeople	0.46	0.77	0.90	1.17	0.96	1.26	0.96	1.23	1.02	1.26	1.04	1.28	
	Johnny	0.39	0.37	0.72	0.69	0.70	0.55	0.81	0.81	0.81	0.78	0.80	0.84	
Average	KristenAndSara	0.40	0.55	0.87	0.86	0.94	0.91	0.92	0.92	0.95	0.96	1.01	0.93	
	Average	0.20	0.56	0.69	1.36	0.74	1.41	0.76	1.48	0.79	1.52	0.82	1.56	
32	Average	0.15	0.34	0.71	0.99	0.76	1.03	0.78	1.05	0.80	1.07	0.83	1.10	
27	Average	0.16	0.23	0.67	0.63	0.69	0.64	0.72	0.67	0.74	0.70	0.77	0.75	
22	Average	0.20	0.13	0.57	0.32	0.56	0.34	0.58	0.35	0.60	0.37	0.62	0.39	

Table 3. Quantitative comparison in terms of Δ PSNR (dB) and Δ SSIM ($\times 10^{-2}$) under RA configuration. **Red** text indicates the best and **blue** text indicates the second best performance.

QP	MFQE	STDF-R3L	RFDA	CF-STIF	STDR	Ours
	PSNR/SSIM	PSNR/SSIM	PSNR/SSIM	PSNR/SSIM	PSNR/SSIM	PSNR/SSIM
37	0.11/0.41	0.43/0.78	0.48/0.90	0.53/1.02	0.55/1.04	0.57/1.07
32	0.12/0.29	0.44/0.55	0.46/0.58	0.49/0.60	0.50/0.63	0.55/0.67
27	0.10/0.20	0.41/0.31	0.42/0.33	0.45/0.40	0.48/0.41	0.53/0.45
22	0.14/0.12	0.33/0.20	0.36/0.23	0.42/0.24	0.46/0.26	0.48/0.27

Table 4. The comparison of parameters and FPS.

Method	Param.(K)	FPS @ Different Resolution		
		832 \times 480	416 \times 240	1280 \times 720
STDF-R3L	1275	11.21	43.53	4.78
RFDA	1270	17.91	59.69	7.08
CF-STIF	2200	4.17	16.32	1.78
STDR	1324	8.77	33.95	3.65
Ours	1386	10.42	39.11	4.08

Table 5. The ablation study of coding priors.

MV P-frame R-frame			Param. (K)	Δ PSNR (dB)	Δ SSIM($\times 10^{-2}$)
Model-1	-	-	1242	0.69	1.40
Model-2	✓	-	1348	0.72	1.42
Model-3	-	✓	1348	0.73	1.44
Model-4	-	-	1280	0.72	1.42
Model-5	✓	✓	1348	0.77	1.50
Model-6	-	✓	1386	0.79	1.52
Model-7	✓	✓	1386	0.82	1.56

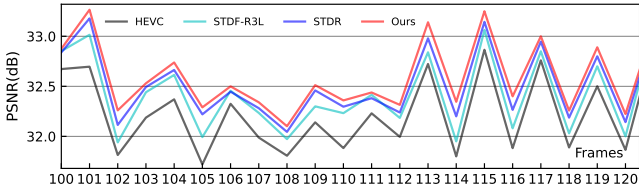


Figure 6. PSNR curves of HEVC, STDF-R3L, STDR and our model on BasketballPass testing sequence under LDB configuration at QP=37.

5.4. Ablation Studies

Effectiveness of coding priors. We firstly implement a baseline model without coding priors (Model-1). Subsequently, we gradually introduce MVs and predictive frames into the IFA module of our baseline and the current residual frame into the MNA module of our baseline to verify the effectiveness of three coding priors (Model-2, Model-3, and Model-4). Besides, we incorporate two coding priors, i.e., MVs and predictive frames or predictive frames and current residual frame to construct two new models, denoting as Model-5 and Model-6, respectively. Their results and corresponding model parameters are provided in Tab. 5. The

average Δ PSNR of all testing sequences under LDB configuration at QP = 37 is used to evaluate the performance gain. It is found from Tab. 5 that our baseline achieves 0.69dB in terms of Δ PSNR. Compared with STDF-R3L [6], our baseline saves 33K parameters and achieves better performance. Furthermore, as can be seen from Tab. 5, MVs, predictive frames and residual frames bring 0.03dB, 0.05dB and 0.05dB performance gain in terms of Δ PSNR, compared to our baseline. Benefiting from incorporating these three coding priors, our model achieves an overall performance gain of 0.13dB in terms of Δ PSNR.

We also provide some feature visual results using different coding priors in Fig. 8. Specifically, we illustrate the temporally-aggregated feature \mathcal{F}^{ta} of Model-1, Model-2 and Model-5, and the spatially-aggregated feature \mathcal{F}^{sa} of

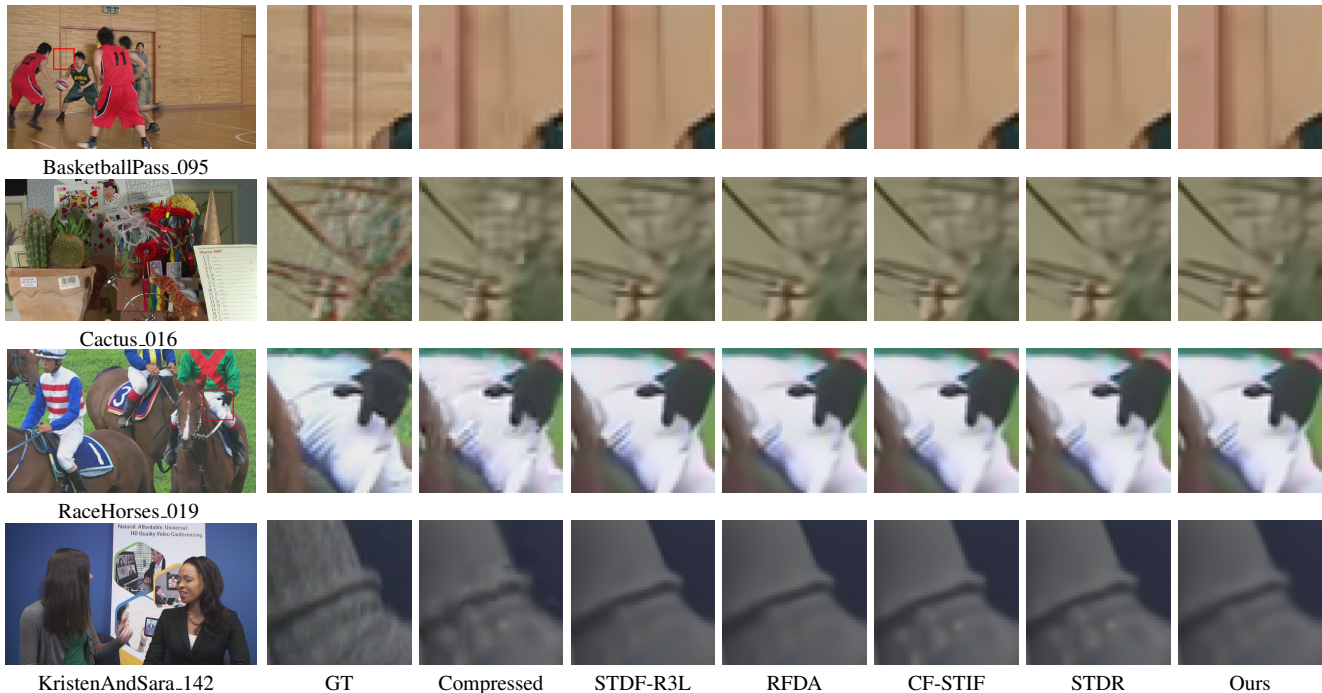


Figure 7. Visual results obtained by using different enhancement methods. The settings in the LQ sequences: BasketballPass at QP = 22 (RA), Cactus at QP = 27 (LDB), RaceHorses QP = 32 (RA), and KristenAndSara at QP = 37 (LDB).

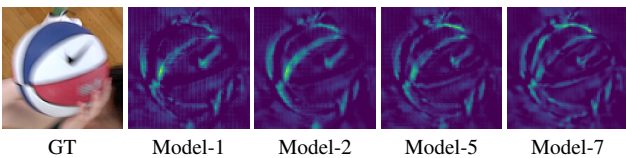


Figure 8. Feature visual results obtained by using coding priors.

Model-7. It is found from Fig. 8 that using MVs and predictive frames, the objects in temporally-aggregated feature \mathcal{F}^{ta} have the more details. Moreover, compared with the results of Model-2 and Model-5, clear details are constructed in the feature results of Model-7, which benefits from introducing the current residual frame and our MNA module.

Effectiveness of VCP dataset. We use the LDP configuration in MFQE dataset [10] to extract three coding priors, i.e., MV, predictive frame and residual frame, from the VCP dataset and denote this dataset as VCP(LDP) dataset. Based on this dataset, we train three state-of-the-art methods, i.e., STDF-R3L, CF-STIF, STDR, and our CPGA model at QP=37 to validate the effectiveness of our VCP dataset under LDP configuration. The corresponding results are given in Tab. 6. In addition, we extract these three coding priors from the MFQE dataset with LDP configuration in [10] under QP=37 and denote this dataset as MFQE-CP to train our CPGA. The corresponding results of our CPGA are provided in Tab. 6. We also offer the results of STDF-R3L [6], CF-STIF [16], and STDR [17] in Tab. 6 on MFQE dataset [10] for fair comparison, which were obtained from their respective papers. It is found from Tab. 6

Table 6. Quantitative comparison on MFQE-CP dataset and VCP(LDP) dataset in terms of Δ PSNR (dB) and Δ SSIM ($\times 10^{-2}$) LDP configuration at QP=37.

Training Dataset	STDF-R3L		CF-STIF		STDR		Ours	
	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
MFQE-CP	0.83	1.53	0.92	1.67	0.98	1.79	1.00	1.82
VCP(LDP)	0.85	1.57	0.95	1.71	1.00	1.81	1.03	1.85

that these compared methods achieve a higher performance gain by using our VCP dataset than the MFQE dataset, which demonstrates the effectiveness of our VCP dataset. It is also shown from Tab. 6 that our CPGA achieves a performance gain of 0.02dB compared to the state-of-the-art method, i.e., STDR, on the MFQE-CP dataset, which further verifies the effectiveness of our CPGA.

6. Conclusion

In this paper, we propose a compressed video quality enhancement dataset that includes LQ videos and their abundant coding priors. Based on this dataset, we design a novel video quality enhancement model, named CPGA. Our proposed CPGA can effectively aggregate inter-frame temporal correlations and spatial correlations with the guidance of these coding priors to generate the HQ videos. Extensive experiments demonstrate that the CPGA achieves the state-of-the-art method for VQE.

Acknowledgement. Thank Shijin Huang for the contributions to the dataset. This work is supported by the National Key Program of China (No. xxx), the National Natural Science Foundation of China (No. xxx).

References

- [1] Christos George Bampis, Zhi Li, Anush Krishna Moorthy, Ioannis Katsavounidis, Anne Aaron, and Alan Conrad Bovik. Study of temporal effects on subjective video quality of experience. *IEEE Transactions on Image Processing*, 26(11):5217–5231, 2017. **1**
- [2] Benjamin Bross, Ye-Kui Wang, Yan Ye, Shan Liu, Jianle Chen, Gary J Sullivan, and Jens-Rainer Ohm. Overview of the versatile video coding (vvc) standard and its applications. *IEEE Transactions on Circuits and Systems for Video Technology*, 31(10):3736–3764, 2021. **1**
- [3] Peilin Chen, Wenhan Yang, Meng Wang, Long Sun, Kangkang Hu, and Shiqi Wang. Compressed domain deep video super-resolution. *IEEE Transactions on Image Processing*, 30:7156–7169, 2021. **2, 3**
- [4] Yue Chen, Debargha Murherjee, Jingning Han, Adrian Grange, Yaowu Xu, Zoe Liu, Sarah Parker, Cheng Chen, Hui Su, Urvang Joshi, et al. An overview of core coding tools in the av1 video codec. In *2018 picture coding symposium (PCS)*, pages 41–45. IEEE, 2018. **1**
- [5] Yuaning Dai, Dong Liu, and Feng Wu. A convolutional neural network approach for post-processing in hevc intra coding. In *MultiMedia Modeling: 23rd International Conference, MMM 2017, Reykjavik, Iceland, January 4-6, 2017, Proceedings, Part I 23*, pages 28–39. Springer, 2017. **1, 2**
- [6] Jianing Deng, Li Wang, Shiliang Pu, and Cheng Zhuo. Spatio-temporal deformable convolution for compressed video quality enhancement. In *Proceedings of the AAAI conference on artificial intelligence*, pages 10696–10703, 2020. **1, 2, 5, 6, 7, 8**
- [7] Chao Dong, Yubin Deng, Chen Change Loy, and Xiaoou Tang. Compression artifacts reduction by a deep convolutional network. In *Proceedings of the IEEE international conference on computer vision*, pages 576–584, 2015. **1, 2**
- [8] Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip Hausser, Caner Hazirbas, Vladimir Golkov, Patrick Van Der Smagt, Daniel Cremers, and Thomas Brox. FlowNet: Learning optical flow with convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2758–2766, 2015. **4**
- [9] Wenbin Du, Yali Wang, and Yu Qiao. Rpan: An end-to-end recurrent pose-attention network for action recognition in videos. In *Proceedings of the IEEE international conference on computer vision*, pages 3725–3734, 2017. **1**
- [10] Zhenyu Guan, Qunliang Xing, Mai Xu, Ren Yang, Tie Liu, and Zulin Wang. Mfqe 2.0: A new approach for multi-frame quality enhancement on compressed video. *IEEE transactions on pattern analysis and machine intelligence*, 43(3): 949–963, 2019. **1, 2, 3, 8**
- [11] Gang He, Kepeng Xu, Chang Wu, Zijia Ma, Xing Wen, and Ming Sun. Hybrid video coding scheme based on vvc and spatio-temporal attention convolution neural network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 1791–1794, 2022. **1**
- [12] Xiaoyi He, Qiang Hu, Xiaoyun Zhang, Chongyang Zhang, Weiyao Lin, and Xintong Han. Enhancing hevc compressed videos with a partition-masked convolutional neural network. In *2018 25th IEEE International Conference on Image Processing (ICIP)*, pages 216–220. IEEE, 2018. **2**
- [13] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. 2014. **3**
- [14] Weiyao Lin, Xiaoyi He, Xintong Han, Dong Liu, John See, Junni Zou, Hongkai Xiong, and Feng Wu. Partition-aware adaptive switching neural networks for post-processing in hevc. *IEEE Transactions on Multimedia*, 22(11):2749–2763, 2019. **2**
- [15] Hongbo Liu, Mingda Wu, Kun Yuan, Ming Sun, Yansong Tang, Chuanchuan Zheng, Xing Wen, and Xiu Li. Ada-dqa: Adaptive diverse quality-aware feature acquisition for video quality assessment. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 6695–6704, 2023. **1**
- [16] Dengyan Luo, Mao Ye, Shuai Li, and Xue Li. Coarse-to-fine spatio-temporal information fusion for compressed video quality enhancement. *IEEE Signal Processing Letters*, 29:543–547, 2022. **1, 2, 6, 8**
- [17] Dengyan Luo, Mao Ye, Shuai Li, Ce Zhu, and Xue Li. Spatio-temporal detail information retrieval for compressed video quality enhancement. *IEEE Transactions on Multimedia*, 2022. **1, 2, 3, 6, 8**
- [18] Di Ma, Fan Zhang, and David R Bull. Bvi-dvc: A training database for deep video compression. *IEEE Transactions on Multimedia*, 24:3847–3858, 2021. **3**
- [19] Jens-Rainer Ohm, Gary J Sullivan, Heiko Schwarz, Thiow Keng Tan, and Thomas Wiegand. Comparison of the coding efficiency of video coding standards—including high efficiency video coding (hevc). *IEEE Transactions on circuits and systems for video technology*, 22(12):1669–1684, 2012. **6**
- [20] Wanli Ouyang, Xiaogang Wang, Xingyu Zeng, Shi Qiu, Ping Luo, Yonglong Tian, Hongsheng Li, Shuo Yang, Zhe Wang, Chen-Change Loy, et al. Deepid-net: Deformable deep convolutional neural networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2403–2412, 2015. **1**
- [21] Woon-Sung Park and Munchurl Kim. Cnn-based in-loop filtering for coding efficiency improvement. In *2016 IEEE 12th Image, Video, and Multidimensional Signal Processing Workshop (IVMSP)*, pages 1–5. IEEE, 2016. **1, 2**
- [22] M. H. Pinson. The consumer digital video library [best of the web]. pages 5646–5654, 2021. **3**
- [23] Rui Qin, Ming Sun, Fangyuan Zhang, Xing Wen, and Bin Wang. Blind image super-resolution with rich texture-aware codebook. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 676–687, 2023. **1**
- [24] Yajun Qiu, Qiang Zhu, Shuyuan Zhu, and Bing Zeng. Dual circle contrastive learning-based blind image super-resolution. *IEEE Transactions on Circuits and Systems for Video Technology*, 34(3):1757–1771, 2024. **1**
- [25] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention—MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18*, pages 234–241. Springer, 2015. **5**

- [26] Kalpana Seshadrinathan, Rajiv Soundararajan, Alan Conrad Bovik, and Lawrence K Cormack. Study of subjective and objective quality assessment of video. *IEEE transactions on Image Processing*, 19(6):1427–1441, 2010. [1](#)
- [27] Gary J Sullivan, Jens-Rainer Ohm, Woo-Jin Han, and Thomas Wiegand. Overview of the high efficiency video coding (hevc) standard. *IEEE Transactions on circuits and systems for video technology*, 22(12):1649–1668, 2012. [1](#), [3](#)
- [28] Thiow Keng Tan, Rajitha Weerakkody, Marta Mrak, Naeem Ramzan, Vittorio Baroncini, Jens-Rainer Ohm, and Gary J Sullivan. Video quality evaluation methodology and verification testing of hevc compression performance. *IEEE Transactions on Circuits and Systems for Video Technology*, 26(1):76–90, 2015. [1](#)
- [29] Dezhao Wang, Sifeng Xia, Wenhan Yang, Yueyu Hu, and Jiaying Liu. Partition tree guided progressive rethinking network for in-loop filtering of hevc. In *2019 IEEE International Conference on Image Processing (ICIP)*, pages 2671–2675. IEEE, 2019. [2](#)
- [30] Tingting Wang, Mingjin Chen, and Hongyang Chao. A novel deep learning-based method of improving coding efficiency from the decoder-end for hevc. In *2017 data compression conference (DCC)*, pages 410–419. IEEE, 2017. [1](#), [2](#), [6](#)
- [31] Yingwei Wang, Takashi Isobe, Xu Jia, Xin Tao, Huchuan Lu, and Yu-Wing Tai. Compression-aware video super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2012–2021, 2023. [2](#), [3](#)
- [32] Thomas Wiegand, Gary J Sullivan, Gisle Bjontegaard, and Ajay Luthra. Overview of the h. 264/avc video coding standard. *IEEE Transactions on circuits and systems for video technology*, 13(7):560–576, 2003. [1](#)
- [33] Mathias Wien. High efficiency video coding. *Coding Tools and specification*, 24, 2015. [1](#)
- [34] Li Xu, Gang He, Jinjia Zhou, Jie Lei, Weiying Xie, Yunsong Li, and Yu-Wing Tai. Transcoded video restoration by temporal spatial auxiliary network. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 2875–2883, 2022. [1](#)
- [35] Xiaozhong Xu, Shan Liu, and Zeqiang Li. Tencent video dataset (tvd): A video dataset for learning-based visual data compression and analysis. *arXiv preprint arXiv:2105.05961*, 2021. [3](#)
- [36] Tianfan Xue, Baian Chen, Jiajun Wu, Donglai Wei, and William T Freeman. Video enhancement with task-oriented flow. *International Journal of Computer Vision*, 127:1106–1125, 2019. [3](#)
- [37] Ren Yang. Ntire 2021 challenge on quality enhancement of compressed video: Dataset and study. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 667–676, 2021. [3](#)
- [38] Ren Yang, Mai Xu, Tie Liu, Zulin Wang, and Zhenyu Guan. Enhancing quality for hevc compressed videos. *IEEE Transactions on Circuits and Systems for Video Technology*, 29(7):2039–2054, 2018. [2](#)
- [39] Ren Yang, Mai Xu, Tie Liu, Zulin Wang, and Zhenyu Guan. Enhancing quality for hevc compressed videos. *IEEE Transactions on Circuits and Systems for Video Technology*, 29(7):2039–2054, 2018. [1](#)
- [40] Ren Yang, Mai Xu, Zulin Wang, and Tianyi Li. Multi-frame quality enhancement for compressed video. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6664–6673, 2018. [1](#), [2](#), [3](#), [6](#)
- [41] Ren Yang, Radu Timofte, Meisong Zheng, Qunliang Xing, Minglang Qiao, Mai Xu, Lai Jiang, Huaida Liu, Ying Chen, Youcheng Ben, et al. Ntire 2022 challenge on super-resolution and quality enhancement of compressed video: Dataset, methods and results. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1221–1238, 2022. [3](#), [6](#)
- [42] Kun Yuan, Zishang Kong, Chuanchuan Zheng, Ming Sun, and Xing Wen. Capturing co-existing distortions in user-generated content for no-reference video quality assessment. In *Proceedings of the 31th ACM International Conference on Multimedia*, pages 1098–1107, 2023. [1](#)
- [43] Bowen Zhang, Limin Wang, Zhe Wang, Yu Qiao, and Hanli Wang. Real-time action recognition with deeply transferred motion vector cnns. *IEEE Transactions on Image Processing*, 27(5):2326–2339, 2018. [2](#)
- [44] Hengsheng Zhang, Xueyi Zou, Jiaming Guo, Youliang Yan, Rong Xie, and Li Song. A codec information assisted framework for efficient compressed video super-resolution. In *European Conference on Computer Vision*, pages 220–235. Springer, 2022. [2](#), [3](#)
- [45] Xiaoming Zhang, Tianrui Li, and Xiaole Zhao. Boosting single image super-resolution via partial channel shifting. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13223–13232, 2023. [5](#)
- [46] Yulun Zhang, Kunpeng Li, Kai Li, Lichen Wang, Bineng Zhong, and Yun Fu. Image super-resolution using very deep residual channel attention networks. In *Proceedings of the European conference on computer vision (ECCV)*, pages 286–301, 2018. [5](#)
- [47] Yulun Zhang, Kunpeng Li, Kai Li, Bineng Zhong, and Yun Fu. Residual non-local attention networks for image restoration. 2019. [2](#)
- [48] Kai Zhao, Kun Yuan, Ming Sun, , and Xing Wen. Zoomvqa: Patches, frames and clips integration for video quality assessment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 1302–1310, 2023. [1](#)
- [49] Kai Zhao, Kun Yuan, Ming Sun, Mading Li, and Xing Wen. Quality-aware pretrained models for blind image quality assessment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22302–22313, 2023. [1](#)
- [50] Liang Zhao, Zhihai He, Wenming Cao, and Debin Zhao. Real-time moving object segmentation and classification from hevc compressed surveillance video. *IEEE Transactions on Circuits and Systems for Video Technology*, 28(6):1346–1357, 2016. [2](#)
- [51] Minyi Zhao, Yi Xu, and Shuigeng Zhou. Recursive fusion and deformable spatiotemporal attention for video compression artifact reduction. In *Proceedings of the 29th ACM inter-*

national conference on multimedia, pages 5646–5654, 2021. [1](#), [3](#), [6](#)

- [52] Ziwei Zheng, Le Yang, Yulin Wang, Miao Zhang, Lijun He, Gao Huang, and Fan Li. Dynamic spatial focus for efficient compressed video action recognition. *IEEE Transactions on Circuits and Systems for Video Technology*, 34(2):695–708, 2024. [1](#)
- [53] Qiang Zhu, Shuyuan Zhu, Guanghui Liu, and Zhenming Peng. Infrared small target detection using local feature-based density peaks searching. *IEEE Geoscience and Remote Sensing Letters*, 19:1–5, 2022. [1](#)
- [54] Qiang Zhu, Pengfei Li, and Qianhui Li. Attention retractable frequency fusion transformer for image super resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pages 1756–1763, 2023. [1](#)
- [55] Xizhou Zhu, Han Hu, Stephen Lin, and Jifeng Dai. Deformable convnets v2: More deformable, better results. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9308–9316, 2019. [2](#), [3](#), [5](#), [6](#)