# Capturing Co-existing Distortions in User-Generated Content for No-reference Video Quality Assessment

Kun Yuan[†]
yuankun03@kuaishou.com
Kuaishou Technology

Zishang Kong[†]
zishang.kong@pku.edu.cn
Peking University

Chuanchuan Zheng
zhengchuanchuan@kuaishou.com
Kuaishou Technology

Ming Sun
sunming03@kuaishou.com
Kuaishou Technology

Xing Wen
wenxing@kuaishou.com
Kuaishou Technology

## ABSTRACT

Video Quality Assessment (VQA), which aims to predict the perceptual quality of a video, has attracted raising attention with the rapid development of streaming media technology, such as Facebook, TikTok, Kwai, and so on. Compared with other sequence-based visual tasks (*e.g.*, action recognition), VQA faces two under-estimated challenges unresolved in User Generated Content (UGC) videos. *First*, it is not rare that several frames containing serious distortions (*e.g.*, blocking, blurriness), can determine the perceptual quality of the whole video, while other sequence-based tasks require more frames of equal importance for representations. *Second*, the perceptual quality of a video exhibits a multi-distortion distribution, due to the differences in the duration and probability of occurrence for various distortions. In order to solve the above challenges, we propose *Visual Quality Transformer (VQT)* to extract quality-related sparse features more efficiently. Methodologically, a Sparse Temporal Attention (STA) is proposed to sample keyframes by analyzing the temporal correlation between frames, which reduces the computational complexity from $O(T^2)$ to $O(T \log T)$. Structurally, a Multi-Pathway Temporal Network (MPTN) utilizes multiple STA modules with different degrees of sparsity in parallel, capturing co-existing distortions in a video. Experimentally, VQT demonstrates superior performance than many *state-of-the-art* methods in three public no-reference VQA datasets. Furthermore, VQT shows better performance in four full-reference VQA datasets against widely-adopted industrial algorithms (*i.e.*, VMAF and AVQT).

## CCS CONCEPTS

• **Computing methodologies** → **Artificial intelligence**; **Computer vision**; **Computer vision tasks**; **Scene understanding**.

## KEYWORDS

video quality assessment, user-generated content, spatiotemporal information, distortions, video Transformer, sparse sampling

## 1 INTRODUCTION

User Generated Content (UGC) has brought evolution to the daily-life consumer domain, which empowers amateurs to become active producers more than consumers. Lower video production cost leads to an explosion of UGC videos on video-sharing platforms, such as FaceBook, Kwai, and so on, which aims to deliver high-quality Quality of Experience (QoE) / Quality of Service (QoS) experience
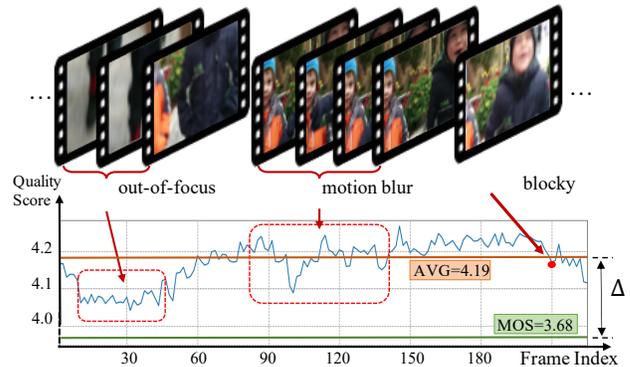


Figure 1: An example of UGC videos of No.13079468475 in the KoNViD-1k [20]. This video contains multiple distortions along the temporal dimension, including motion-related types (*i.e.* out-of-focus and blurriness) and compression-related types (*i.e.* blocking artifacts). There exists a large difference in perceptual quality among frames. And those frames with low quality determine the overall quality instead of an arithmetic mean. Therefore, a more appropriate strategy is needed for spatiotemporal representation in VQA.

to users. Compared with Professionally Generated Content (PGC), UGC videos inevitably have worse conditions of shooting, poor capturing equipment, and unstable transmission links [54]. For video streaming services, VQA has attracted more attention to filter out videos with low perceptual qualities [14, 17]. Furthermore, VQA is also used to conduct content-aware video encoding [6] or enhancement [12, 21, 66], resulting in lower bandwidth cost and better viewing experience. Therefore, it is of great economic value to rate the perceptual quality of UGC videos through VQA.

A large number of studies on image/video quality assessment (QA) are studied in the previous literature. According to the *availability of reference videos*, QA measures can be categorized into full-reference (FR) [16, 26], reduced-reference (RR) [37] and no-reference (NR) [25, 30, 40, 57]. Since distortion-free reference videos are often hard to obtain, NR-VQA is widely adopted in the UGC scenarios. According to the *feature generation types*, QA methods are divided into traditional hand-crafted [38, 40] and learning-based [4, 18, 23, 31] ones. With the rapid development of deep learning,

---

[†]Authors contributed equally to this research.

convolutional neural networks (CNN) [29, 30, 67] and Transformers [25, 57, 59] are also used to boost the VQA domain.

As shown in Fig. 1, the common phenomenon is that multiple distortions co-occur within a UGC video, where different distortions begin to appear at different frames and own different time-span. Such a phenomenon casts two challenges for a better fit of human perceptual quality. **First**, the perceptual quality of a video is determined by the keyframes that contain particular distortions. Excessively dense sampling brings an unbalanced distribution of frames and may disturb the learning process of distortion characteristics. While relatively current sparse sampling may ignore keyframes. How to select frames efficiently is an essential problem to be solved. **Second**, due to the differences in temporal duration of different distortions, the perceptual quality of frames within a video exhibits a multi-distribution mode. Take some distorted characteristics for example, blocking artifact [23], dirty lens [19], and noise [33, 70] are usually easy-detected given an individual frame. But out-of-focus and motion blurriness [34] can only be recognized using multiple frames. These factors put forward a higher request for VQA methods with the ability to perform frames analysis under different durations simultaneously.

To overcome the aforementioned annoying challenges, we propose *Visual Quality Transformer (VQT)* to extract quality-aware features focusing on multi-distortions more efficiently. Specifically, **to solve the first challenge**, a *Sparse Temporal Attention (STA)* is proposed to sample keyframes via analyzing the temporal correlation between frames. It reformulates self-attention from the perspective of sparse sampling and adopts a proper sampling ratio according to the Johnson-Lindenstrauss (JL) lemma [22]. The keyframes can be selected by comparing the Kullback–Leibler (KL) difference between the Uniform distribution and its cosine similarity with other frames. As for model efficiency, compared with vanilla temporal attention, STA reduces the computational complexity from $O(T^2)$ to $O(T \log T)$. **To solve the second challenge**, owing to the efficiency of the STA module, *Multi-Pathway Temporal Network (MPTN)* is adopted to capture co-existing distortions in a video simultaneously, which stacks multiple STA modules with different degrees of sparsity. Finally, the aggregated features are used for the representation of a video in VQA.

Our **contributions** are summarized as follows:

- We propose an effective and efficient Visual Quality Transformer (VQT) for the NR-VQA tasks, where the proposed STA selects key frames containing particular distortions and the MPTN helps capture different distorted characteristics simultaneously in UGC scenarios.
- VQT demonstrates superior performance than many *state-of-the-art (SoTA)* methods in three NR-VQA datasets, raising the performance by 2.14% of PLCC in KoNViD-1k and 2.17% of PLCC in YouTube-UGC over the best results. Furthermore, VQT obtains better results in four FR-VQA datasets (cross-dataset evaluation in three of them) against widely-adopted industrial algorithms (*e.g.*, VMAF and AVQT).
- VQT can act as a plug-in module used for general computer-vision tasks and shows good generalization ability in the video classification task. Compared with the original dense

attention mechanism (*e.g.*, TimeSformer), the computational cost decreases from 197 TFLOPs to 154 TFLOPs (-22%).

## 2 RELATED WORK

### 2.1 Perceptual Quality Assessment

According to the accessibility of the reference images or videos, QA is divided into FR-QA, RR-QA, and NR-QA tasks. FR-QA and RR-QA tasks require a full and partial reference respectively. While the NR-QA method only takes distorted images or videos as input, which is often more challenging, but also more practical in most scenarios. In this paper, we focus on the NR-VQA domain.

In the early period of NR-VQA, most works [1, 36, 56] focused on identifying specific types of distortions(*e.g.*, blur, blocky). Then, more methods [38, 62] have been proposed to focus on multiple distortions jointly to carry out comprehensive QA. With the rapid progress of deep learning, learning-based methods [27, 49, 63, 67–69] have suppressed the performance of traditional hand-crafted ones, due to their versatility and generalization. RAPIQUE [49] combined quality-aware features of scene statistics and semantics-aware deep convolutional features. A combination of 3D-CNN and LSTM was adopted [64] to extract local spatiotemporal features from small clips in the video. Patch-VQA [63] devised a local-to-global patch-based architecture, and extracted both 2D and 3D video features using a temporal DNN to predict the quality. STDAM [60] introduced using the graph convolution and attention module to extract and enhance the quality-related features. 2BiVQA [44] proposed using two Bi-directional Long Short Term Memory (Bi-LSTM) to conduct the quality assessment. One is for capturing short-range dependencies between image patches, and the other is for capturing long-range dependencies between frames. However, how to efficiently select keyframes to model quality-related features is still an open problem in the VQA domain. Recently, DisCoVQA [57] designs a Transformer-based Spatial-Temporal Distortion Extraction module to tackle temporal quality attention.

### 2.2 Video Transformer Architecture

Visual Transformers [10] are the most popular alternatives applied in various vision tasks, due to their good ability in modeling long-term dependency of sequential data. Following ViT, many Transformer-based models [2, 3, 13, 35, 65] were developed for video classification tasks. TimeSformer [3] explored the factorization of spatial-temporal dimension for efficient computation. Video Swin Transformer [35] globally connected patches across the spatial and temporal dimensions, and advocated an inductive bias of locality for a better speed-accuracy trade-off. There also exists some work exploring the applicability of Transformers in the field of VQA. B-VQA [28] combined GRU unit with Transformer encoder to model the temporal information, which further boosts the performance. StarVQA [59] transferred the divided space-time attention in TimeSformer directly into VQA, showing well generalization ability in the regression task. But these architectures are not designed specifically for VQA, especially for efficiently modeling co-existing distortions. In this paper, we deeply analyze the problems faced by VQA and design the VQT in a targeted manner.
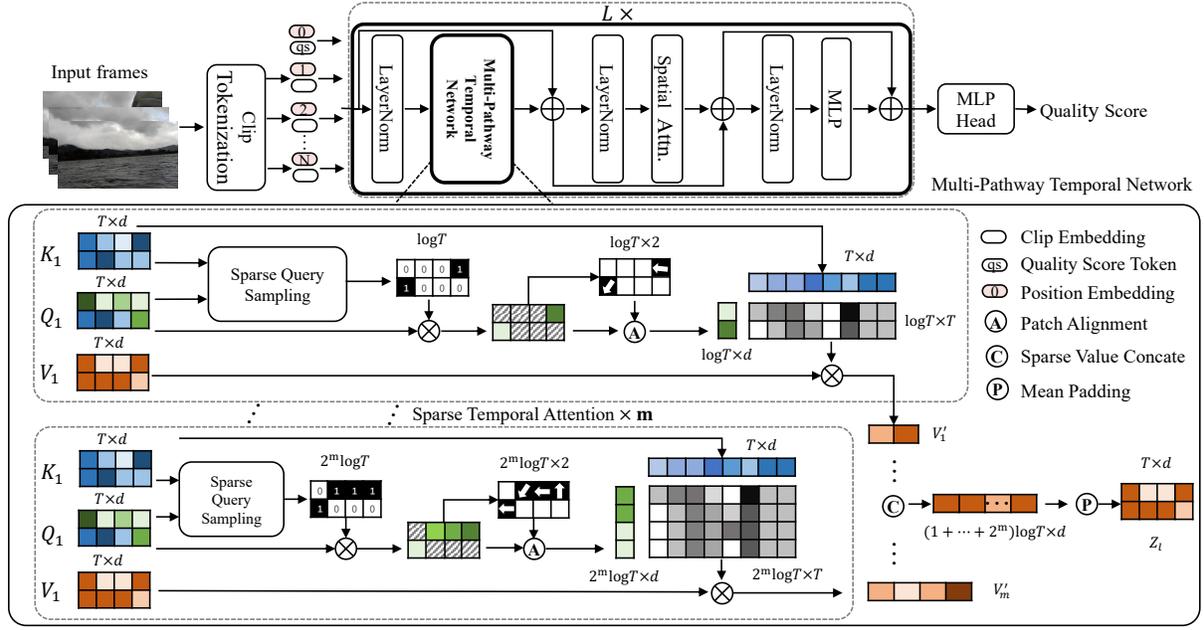
**Figure 2: Illustration of the proposed Visual Quality Transformer (VQT). It receives a clip as input and reshapes each frame into patches for tokenization. Then the sequence of tokens is fed into a stacked encoder, performing the spatiotemporal attention. In the temporal dimension, VQT utilizes a Multi-Pathway Temporal Network to capture different distorted characteristics simultaneously with stacking Sparse Temporal Attention blocks. In each block, STA conducts sparse query sampling to locate keyframes for distortion representation. To further enhance STA with spatial relationships, temporal offsets across frames are assigned to patches for alignment. Different blocks contribute to the final representation within an encoder block. To predict the video quality, the corresponding quality score token is used for the final representation.**

## 3 METHOD

### 3.1 Revisiting Video Transformer

Different from image Transformers, video Transformers receive a sequence of frames as input. There exist some differences in basic architecture modules, including clip tokenization and spatiotemporal attention. We first briefly revisit them for a better understanding.

*Clip Tokenization.* The video Transformer takes a clip as input denoted as $\mathbf{X} \in \mathbb{R}^{T \times H \times W \times 3}$, which composes of $T$ RGB frames with the size of $H \times W$ under equal temporal interval sampled from original videos. Following ViT, video Transformers reshape each frame into $N$ non-overlapping patches, where each size is $P \times P$ and $N = \frac{H \times W}{P^2}$. Then the sequence of frames can be flattened into $\mathbf{X} \in \mathbb{R}^{T \times N \times (3P^2)}$. Besides, an extra learnable positional embedding $\mathbf{E}^{pos}$ is added to encode the spatiotemporal position of each patch. Then the input embedding $\mathbf{Z} \in \mathbb{R}^{T \times N \times d}$ is calculated as:

$$\mathbf{Z} = \mathbf{W}\mathbf{X}^\top + \mathbf{E}^{pos}, \tag{1}$$

where $\mathbf{W} \in \mathbb{R}^{d \times 3P^2}$ is the mapping weight and $d$ is the embedding dimension of each frame token.

*Divided Space-Time Attention.* To process spatiotemporal information, TimeSformer utilizes the "Divided Space-Time Attention" module, where the temporal attention and the spatial attention are performed sequentially. In each encoder block $\ell$, the temporal attention first computes the relationship among all patches in the same spatial location from different frames, expressed as:

$$\mathbf{Z}^\ell_{\text{time}} = \text{Softmax}\left(\frac{\mathbf{Q}^{\ell-1}_{t'}}{\sqrt{d}}\mathbf{K}^{\ell-1\top}_{t'}\right)\mathbf{V}^{\ell-1}, \tag{2}$$

where $\mathbf{Q}$, $\mathbf{K}$, $\mathbf{V}$ are the query, key, and value. And $t'$ denotes that the inner products are computed on the temporal dimension. Then the generated $\mathbf{Z}^\ell_{\text{time}}$ is fed back for the spatial attention, computing as:

$$\mathbf{Z}^\ell \text{space} = \text{Softmax}\left(\frac{\mathbf{Q}^{\ell-1}_{p'}}{\sqrt{d}}\mathbf{K}^{\ell-1\top}_{p'}\right)\mathbf{V}^{\ell-1}, \tag{3}$$

where $p'$ denotes inner products on the spatial dimension.

### 3.2 Visual Quality Transformer

The illustration of VQT is given in Fig. 2. The core components are STA and MPTN, which we will describe in more detail.

*Sparse Temporal Attention.* Video clips own widespread information redundancy in the temporal dimension [55, 61], both in the frame level and feature level. Frames containing distorted characteristics largely reflect the perceptual quality of the whole video. The upper part of Fig. 3 demonstrates this phenomenon in the VQA domain, where the computed temporal attention map shows that all frames have a strong correlation with the 5-*th* frame (*i.e.*, , an over-exposed distorted image). And the factor of overexposure in
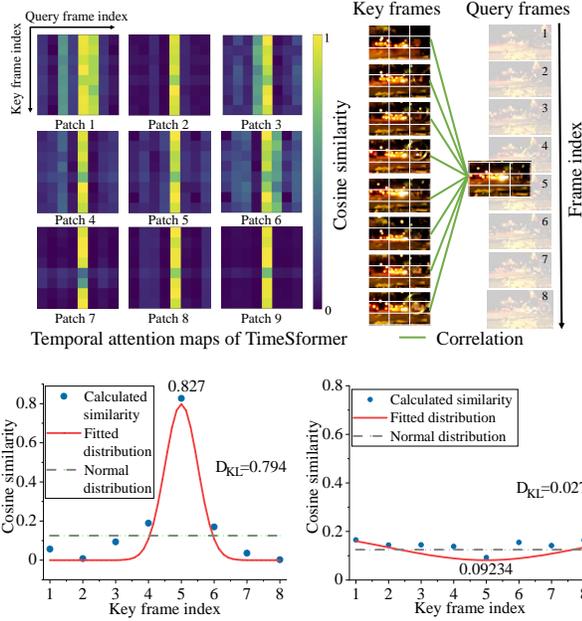
Figure 3: The upper images show that the quality of the video is largely affected by frames containing distortions. The below curves give the distributions of $D_{KL}$ for the 5-*th* (left) and 3-*th* (right) frame.

this frame deteriorates the quality of this video. Based on the observations in Fig. 1 and Fig. 3, we propose a Sparse Temporal Attention, aiming to sample key-frames containing distortions in a video.

Specifically, the attention mechanism is first reformatted from the perspective of sparse sampling. Given the query $\mathbf{Q} \in \mathbb{R}^{T \times d}$, key $\mathbf{K} \in \mathbb{R}^{T \times d}$, and value $\mathbf{V} \in \mathbb{R}^{T \times d}$, the self-attention is computed as $\mathbf{V}' = \text{softmax}(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d}})\mathbf{V}$, where $\mathbf{V}' \in \mathbb{R}^{T \times d}$. According to the JL lemma, *there exists a linear transformation that projects* $\mathbf{V}$ *into* $\mathbf{V}' \in \mathbb{R}^{\log T \times d}$ *with low-distortion embedding*, where $\log T$ is the minimal number of frames sampled from $T$ frames. Specifically, we let $q_i$ and $k_j$ represent the feature of $i$-th row and $j$-th row in $\mathbf{Q}$ and $\mathbf{K}$, respectively. Then we calculate the cosine similarity between $q_i$ and $k_{j,j=1,..,T}$. And the computed distribution can be noted as $p(\mathbf{K}|q_i)$, representing the correlation of the $i$-th frame with all frames. Then we compute the KL divergence between $p(\mathbf{K}|q_i)$ and the uniform distribution $U$:

$$D_{\text{KL}}(p(\mathbf{K}|q_i)\|U) = \sum_{j=1}^{T} \frac{1}{T} \log \frac{1}{T} - \frac{1}{T} \log \left( \frac{q_i k_j^T}{\sum_{j=1}^{T} q_i k_j^T} \right) \quad (4)$$

$$= -\log \frac{1}{T} - \sum_{i=1}^{T} \frac{1}{T} \left( \ln \exp \frac{q_i k_j^T}{\sqrt{d}} - \ln \sum_{j=1}^{T} e^{\frac{q_i k_j^T}{\sqrt{d}}} \right). \quad (5)$$

For simplicity, we ignore the constant term of $-\log \frac{1}{T}$. And the equation can be rewritten as:

$$D_{KL}(p(\mathbf{K}|q_i)\|U) = \ln \sum_{j=1}^{T} e^{\frac{q_i k_j^\top}{\sqrt{d}}} - \frac{1}{T} \sum_{j=1}^{T} \left( \frac{q_i k_j^\top}{\sqrt{d}} \right). \quad (6)$$

**Algorithm 1** Pseudo-code of selection of keyframes

---
1: $Q_{sample} = random.sample(Q, gamma)$ ▷ gamma:logT/T
2: $mu = KL\text{-}Divergence(Q_{sample}, K).mean(dim = -1)$
3: $sigma = KL\text{-}Divergence(Q_{sample}, K).std(dim = -1)$
4: init $i = 0$; $attn = zeros(logT, T)$
5: **for** $q$ in $Q$ **do**
6:     **if** $i < logT$ and $KL\text{-}Divergence(q, K) > mu + sigma$ **then**
7:         $attn[i,:] = q@K.transpose()$
8:         $i += 1$
9:     **end if**
10: **end for**
11: $V = bmm(attn, V)$ ▷ bmm:batch matrix multiplication

---

Larger values of $D_{KL}(p(\mathbf{K}|q_i)\|U)$ indicate a **stronger correlation with other frames**, which can be used for the selection of keyframes. Examples are given in Fig. 3, where frames containing distortions own larger (*e.g.*, , the 5-*th* frame) values and vice versa.

To reduce the complexity of computing the divergence of all frames, $\log T$ frames are sampled randomly. The distribution of divergence of all frames can be estimated by calculating the mean $\mu$ and variance $\sigma$ of $\log T$ frames. Then $\log T$ query frames can be selected, whose divergence meets $D_{KL}(p(\mathbf{K}|q_i)\|U) > \mu + \sigma$. The dimension of sampled query $\hat{\mathbf{Q}}$ is reduced to $\log T \times d$. And the pseudocode is shown in Alg. 1.

To enhance the spatial relationship of queries between different frames, STA further performs a spatial shift for alignment. Specifically, we reshape the spatial dimension of selected query features $\hat{\mathbf{Q}} \in \mathbb{R}^{\log T \times N \times d}$ into $\mathbb{R}^{\log T \times \sqrt{N} \times \sqrt{N} \times d}$ (the spatial dimension $N$ is not shown in Fig. 2 for simplify). Then a linear projection layer is attached on the reshaped features to predict 2D offsets $\mathbf{P}$ for each token, where $\mathbf{P} \in \mathbb{R}^{\log T \times \sqrt{N} \times \sqrt{N} \times 2}$. Then the shifted query features $\hat{\mathbf{Q}}'$ can be obtained by bilinear interpolation. Then the weighted value features can be computed by:

$$\mathbf{V}' = \text{softmax}(\frac{\hat{\mathbf{Q}}'\mathbf{K}^\top}{\sqrt{d}})\mathbf{V}. \quad (7)$$

*Multi-Pathway Temporal Network.* To capture different distorted characteristics simultaneously, multiple STA modules with different degrees of sparsity are stacked in parallel. Given a video clip with $T$ frames, the number $m$ of pathway is determined by $m = \lfloor \log(\frac{T}{\lceil \log T \rceil} + 1) \rfloor - 1$, where $\lceil \cdot \rceil$ and $\lfloor \cdot \rfloor$ represents the ceiling and flooring operation respectively. The minimal number of frames for different STA modules is $\log T$, and the maximal is $2^m \log T$. Take an input clip with 96 frames, for example, its $m$ is 3, containing 3 pathways. And each pathway selects 7, 14, and 28 keyframes respectively. Each block performs temporal attention over selected frames, resulting in a weighted value of $\mathbf{V}'_m \in \mathbb{R}^{2^m \log T \times d}$. Then the values generated by different blocks are concatenated in the temporal dimension, resulting in $\hat{\mathbf{V}} \in \mathbb{R}^{(1+\cdots+2^m) \log T \times d}$. To align with succeeding encoder blocks, we use the mean padding operation to fill $\hat{\mathbf{V}}$ in the temporal dimension, generating $\mathbf{Z}_l \in \mathbb{R}^{T \times d}$:

$$\hat{\mathbf{V}} = \text{Concate}(\mathbf{V}'_1, \cdots, \mathbf{V}'_m)$$
$$\mathbf{Z}_l = \text{Mean-Padding}(\hat{\mathbf{V}}). \quad (8)$$

## 3.3 Optimization Objective

A smooth $\mathcal{L}_1$ loss is adopted to train VQT models. Let $\mathcal{F}(\cdot)$ represent the mapping function of the VQT model. Given mini-batch videos during training, the objective function can be denoted as:

$$\min \ \frac{1}{|\mathcal{B}|} \sum_{i=1}^{\mathcal{B}} \mathcal{L}_{1-smooth}(\mathcal{F}(\mathbf{X}_i), y_i), \qquad (9)$$

where $\mathbf{X}_i$ and $y_i$ is the input video and corresponding labeled Mean Opinion Score (MOS). $\mathcal{B}$ indicates the size of the mini-batch.

## 3.4 Computational Efficiency

*Computational Complexity of STA.* For an input clip with $T$ frames, the computational cost of the original temporal attention module is $O(T^2)$. The STA module applies sparse computation among the temporal dimension with a cost of $O(2^m \log T \cdot T)$. And $2^m \log T$ is smaller than $T$ as mentioned above.

*Computational Complexity of MPTN.* MPTN is composed of multiple STA blocks, whose total computational cost is computed by combining each one of $O(\sum_{a=0}^{m} 2^a \log T \cdot T \cdot d)$. Since $\sum_{a=0}^{m} 2^a \log T$ is smaller than $T$, the computational cost of MPTN is still less than that of the original temporal attention module.

*Measured Inference Speed.* Further evaluations of efficiency are given in the following experiments. In Tab. 5 and 6, compared with the original version of dense attention (*i.e.*, , TimeSformer used in StarVQA), VQT has a less computational cost (-22%), faster inference speed (+13%) and higher performance in PLCC (+3.49%).

## 4 EXPERIMENTS

### 4.1 Datasets and Evaluation

*NR datasets.* We leverage three NR-VQA datasets to evaluate VQT models: LIVE Video Quality Challenge Database (LIVE-VQC) [43], Konstanz Natural Video Database (KoNViD-1k) [20], and Blind Video Quality Assessment for User Generated Content (Youtube-UGC) [53]. Subjective quality scores are provided in the form of MOS. LIVE-VQC contains 585 videos labeled by MOS with a resolution from 240P to 1080P. KoNViD-1k comprises a total of 1,200 videos with a resolution of $960 \times 540$ that are fairly sampled from a large public video dataset, YFCC100M[45]. The duration of videos is 8s with 24/25/30FPS, whose MOS ranges from 1.22 to 4.64. Youtube-UGC is composed of 1,500 videos that are sampled from millions of YouTube videos belonging to 15 categories annotated by a knowledge graph. The resolutions of videos are from 360P to 4K. For the Youtube-UGC, we follow the default training and testing splits [53]. For LIVE-VQC and KoNViD-1k, following [49], 80% of the dataset is used for training, and the remaining 20% is used for testing.

*FR Datasets.* To verify the generalization ability to video codec field, four more FR-VQA datasets are tested, including ICME-FR [†], VQEG HD3 [†], NFLX Video dataset [†] and the Waterloo IVC 4k Video Quality database [32]. The VQEG HD3 database is composed of 9 source clips with a resolution of 1080P. And the source clips are encoded into 63 distorted videos for evaluation. The NFLX Video

[†]http://2021.ieeeicme.org/2021.ieeeicme.org
[†]https://www.cdvl.org
[†]https://github.com/Netflix/vmaf

consists of 34 source clips, whose duration is 6s. They are sampled from popular TV shows and movies on Netflix. Source clips are encoded at resolutions ranging from $384 \times 288$ to $1920 \times 1080$, resulting in about 300 distorted videos. The Waterloo IVC 4k Video Quality database is created from 20 pristine 4K videos. Each video is encoded by five encoders: HEVC, H264, VP9, AV1, and AVS2, and divided into three solutions ($960 \times 540$, $1920 \times 1080$ and $3840 \times 2160$ with four distortion levels, resulting in 1,200 encoded videos.

*Evaluation Criteria.* Pearson's Linear Correlation Coefficient (PLCC), Spearman's Rank-Order Correlation Coefficient (SROCC), Kendall's Rank-Order Correlation Coefficient (KROCC), and Root Mean Square Error (RMSE) are used for evaluation. PLCC and RMSE measure the prediction accuracy, SROCC, and KROCC indicate the prediction monotonicity. Better VQA methods should have larger PLCC/SROCC/KROCC and smaller RMSE values.

### 4.2 Implementation Details

Our implementation is based on PyTorch [41] and MMAction2 [8]. All models are trained using 4 NVIDIA Tesla V100. The number of encoders follows the original setting of TimeSformer. We set a patch size of 16 in clip tokenization. The embedding dimension $d$ is 768. We use models that have been pre-trained on ImageNet [9] and Kinetics-400 [24] for training. During the optimization procedure, we use the AdamW optimizer with a learning rate of 1e-5 decayed by a factor of 0.1 every 30 epochs, minimizing the $\mathcal{L}1$ loss. All models are trained for 90 epochs. By default, the checkpoint generated by the last iteration is used for evaluation. The batch size of video clips is set to 4 with a clip length of 96. Other training settings are the same with [3]. The median result of 10 repeat runs is used for Tab. 1 with different random splits.

### 4.3 Comparison with SoTA Methods

Extensive experiments are conducted to compare with SoTA QA methods. As given in Tab. 1, we report the PLCC and SROCC performance in KoNViD-1k, LIVE-VQC and Youtube-UGC datasets. Besides, weighted average scores are reported based on the number of videos of three datasets. Some observations and conclusions can be found here. *First*, compared with IQA methods (i.e. NIQE [40], BRISQUE [38], CORNIA [62]), VQT obtains a large performance lead (+31.54%, +24.24%, +26.04% of PLCC in KoNViD-1k), showing the effectiveness of fusion strategy of frames over an arithmetic mean. *Second*, compared with hand-crafted features (i.e. [39]), VQT demonstrates the utility of learning-based methods over prior knowledge. *Third*, compared with CNN models (i.e. [28, 49, 60, 63]), VQT shows the advantage of Transformer models in building long-range dependencies. *Fourth*, compared with current Transformer models (i.e. StarVQA [59], TimeSformer), VQT also evaluate the gains from architecture modification by large margins (+7.24% of PLCC in KoNViD-1k, +2.77% of PLCC in LIVE-VQC). Compared with current SOTA method, VQT also outperforms STDAM [60] in three datasets(+2.59% of PLCC in KoNViD-1k, +1.53% of PLCC in LIVE-VQC, and +2.17% of PLCC in Youtube-UGC). Since the size of LIVE-VQC is the smallest (only containing 585 videos), BVQA [28] achieves the highest performance by introducing extra QA training data. However, VQT still surpasses it in the weighted scores (+2.39%

**Table 1: Quantitative results of different methods on three public NR-VQA datasets. Larger PLCC and SROCC indicate better performance. Besides, weighted average scores are reported based on the number of videos of three datasets. The best and second best performances are highlighted and underlined. The mark "-" denotes that results are not reported originally. The "*" mark indicates using extra training data for QA tasks. The VQT models outperform almost all SoTA methods by large margins.**

| Method | KoNViD-1k | | LIVE-VQC | | Youtube-UGC | | Weighted Average | |
|---|---|---|---|---|---|---|---|---|
| | PLCC ↑ | SROCC ↑ | PLCC ↑ | SROCC ↑ | PLCC ↑ | SROCC ↑ | PLCC ↑ | SROCC ↑ |
| VIIDEO [39] | 0.303 | 0.298 | 0.2164 | 0.0332 | 0.1534 | 0.0580 | 0.2230 | 0.1459 |
| NIQE [40] | 0.5530 | 0.5417 | 0.6286 | 0.5957 | 0.2776 | 0.2379 | 0.4500 | 0.4225 |
| CORNIA [62] | 0.608 | 0.610 | - | - | - | - | - | - |
| BRISQUE [38] | 0.626 | 0.654 | 0.638 | 0.592 | 0.395 | 0.382 | 0.5299 | 0.5265 |
| VBLIINDS [42] | 0.6576 | 0.6947 | 0.7120 | 0.7015 | 0.5551 | 0.5590 | 0.6431 | 0.6427 |
| GRU-VQA [29] | 0.744 | 0.755 | - | - | - | - | - | - |
| TLVQM [11] | 0.7688 | 0.7729 | 0.8025 | 0.7988 | 0.6590 | 0.6693 | 0.7284 | 0.7337 |
| MDTVSFA [30] | 0.7856 | 0.7812 | 0.7728 | 0.7382 | - | - | - | - |
| UGC-VQA [48] | 0.7803 | 0.7832 | 0.7514 | 0.7522 | 0.7733 | 0.7787 | 0.7719 | 0.7754 |
| PVQ [63] | 0.786 | 0.791 | 0.837 | 0.827 | - | - | - | - |
| RAPIQUE [49] | 0.8175 | 0.8031 | 0.7863 | 0.7548 | 0.7684 | 0.7591 | 0.7907 | 0.7753 |
| StarVQA [59] | 0.796 | 0.812 | 0.808 | 0.732 | - | - | - | - |
| BVQA* [28] | 0.8335 | 0.8362 | **0.8415** | **0.8412** | 0.8194 | 0.8312 | 0.8290 | <u>0.8350</u> |
| STDAM [60] | <u>0.8415</u> | <u>0.8448</u> | 0.8204 | 0.7931 | <u>0.8297</u> | <u>0.8341</u> | <u>0.8320</u> | 0.8305 |
| 2BiVQA [44] | 0.835 | 0.815 | 0.832 | 0.761 | 0.790 | 0.771 | 0.794 | 0.800 |
| DisCoVQA [57] | 0.847 | 0.847 | 0.826 | 0.820 | - | - | - | - |
| Our TimeSformer | 0.8293 | 0.8342 | 0.8017 | 0.7845 | 0.8279 | 0.8133 | 0.8235 | 0.8159 |
| VQT | **0.8684** | **0.8582** | <u>0.8357</u> | <u>0.8238</u> | **0.8514** | **0.8357** | **0.8529** | **0.8421** |

of PLCC), showing strong generalization ability. Compared with recent DisCoVQA, VQT obtain a higher result of PLCC by 2.14%.

## 4.4 Comparing with VMAF/AVQT

We further compare the effectiveness of our method with the two industrial standards, i.e. Netflix's VMAF [†] and Apple's AVQT [†] on the four widely adopted open datasets. The VMAF and AVQT have been used as standard deals to their simplicity in computations and consistent performance across different types of videos. As shown in Tab. 2, their performance is consistently high across the four datasets except for the *Waterloo IVC 4k* since both VMAF and AVQT were developed before 4K videos were becoming popular. The performance decrease on *Waterloo IVC 4k* indicates that the generability of these two algorithms is not always satisfying when a new video format is induced.

We train our VQT model with data solely from the ICME-FR train split. However, we not only test it on the ICME-FR test split but also directly evaluate its performance on the aforementioned three datasets without any fine-tuning (*i.e.*, **cross-dataset evaluation**). From Tab. 2, we observe that our proposed VQT method further improves the performance. This enhancement is not only attributed to the fact that our model is learned from a large-scale dataset but also thanks to our carefully designed architecture that can effectively extract and integrate spatiotemporal features to better model users' Mean Opinion Scores (MOS). It is worth noting that our method is efficient and can be easily integrated and substituted by existing

**Table 2: Comparisons with industrial standards. The VQT model is trained on the ICME-FR datasets with the supervision of DMOS. And direct inference results are reported on the other three datasets without fine-tuning. VQT shows strong generalization ability and practical prospects.**

| Datasets | Method | PLCC ↑ | SROCC ↑ |
|---|---|---|---|
| ICME-FR | VMAF | 0.9423 | 0.9137 |
| | AVQT | 0.9730 | 0.9334 |
| | VQT | **0.9867** | **0.9364** |
| NFLX Video Dataset | VMAF | 0.9351 | 0.9173 |
| | AVQT | 0.9571 | 0.9420 |
| | VQT | **0.9715** | **0.9532** |
| VQEG HD3 | VMAF | 0.9266 | 0.9238 |
| | AVQT | 0.9481 | 0.9417 |
| | VQT | **0.9603** | **0.9576** |
| Waterloo IVC 4k | VMAF | 0.7324 | 0.7325 |
| | AVQT | 0.7749 | 0.7738 |
| | VQT | **0.7885** | **0.7821** |

methods such as VMAF/AVQT for evaluating the performance of different video encoding strategies.

## 4.5 Ablation Studies and Visualization

To verify the rationality of the proposed modules, ablation studies are conducted in the following aspects.

---

[†]https://github.com/Netflix/vmaf
[†]https://developer.apple.com/videos/play/wwdc2021/10145

**Table 3: Ablation study of individual modules in VQT, conducted in KoNViD-1k, Youtube-UGC, and Kinetics-400.**

| Modules | | KoNViD-1k | | YoutuUGC | | K400 |
|---|---|---|---|---|---|---|
| STA | MPTN | PLCC ↑ | SROCC ↑ | PLCC ↑ | SROCC ↑ | Top-1 Acc ↑ |
| ✓ | ✓ | **0.8684** | **0.8582** | **0.8514** | **0.8357** | **80.3** |
| ✓ | ✗ | 0.8530 | 0.8510 | 0.8429 | 0.8280 | 79.0 |
| ✗ | ✗ | 0.8293 | 0.8342 | 0.8279 | 0.8133 | 78.0 |

**Table 4: Ablation study on the type of reduction and the number of frames used in the STA module. Experiments are conducted in the KoNViD-1k.**

| Sampling Strategy | Frames | PLCC↑ | SROCC↑ |
|---|---|---|---|
| our KL-based | $\log T$ | **0.8684** | **0.8582** |
| our KL-based | $0.5 \log T$ | 0.8320 | 0.8309 |
| our KL-based | $2 \log T$ | 0.8691 | 0.8575 |
| Random[7] | $\log T$ | 0.8067 | 0.7798 |
| Linear[50] | $\log T$ | 0.8142 | 0.8091 |
| Conv[51] | $\log T$ | 0.8140 | 0.8090 |
| Clustering | $\log T$ | 0.8379 | 0.8203 |

*The rationale of the proposed STA module.* The rationale of STA is verified *theoretically and experimentally. Theoretically,* according to the JL lemma, the *lower bound* of the error coefficient $\varepsilon$ after projection is measured by $d > 8 \ln(T)/\varepsilon^2$, where $T$ is the number of frames, and $d$ is the embedding size. In our setting of $T = 96$ and $d = 768$. So $\varepsilon$ can be calculated as 0.215, which means **more than 78.5%** of temporal information or more is maintained by selected $\log T$ keyframes. *Experimentally,* sufficient ablation studies also confirm the effectiveness of STA as shown in Tab. 3 and 4. Simply adding the STA module to the baseline can bring consistent promotion on two VQA datasets (KoNViD-1k and Youtube-UGC) and one general video classification dataset (Kinetics-400).

*Reduction types in the STA module.* We conducted an evaluation of various linear reduction methods, including *random*[7] ($\log T$ frames are randomly selected from $T$ frames), *linear reduction*[50] (features are transformed using a matrix of $\log T \times T$, resulting in the representation of $\log T$ frames), *conv reduction*[51] (features are transformed using a Conv/BN/ReLU module, which reduces the channel from $T$ to $\log T$), *clustering* (features are clustered into $\log T$ centers according to cosine similarity), and *STA*. The results, presented in Tab. 4, demonstrate that the *STA* reduction method has the most significant positive impact on performance, effectively removing redundant information. Additionally, we found that the $\log T$ setting was the optimal sparsity setting through experimentation involving increasing and decreasing the number of frames used to represent keyframes.

*Different combinations of proposed modules.* We conduct ablation experiments for different proposed modules, including STA and MPTN. Results are given in Tab. 3. The best performance is

**Table 5: Comparison of inference cost with academic and industrial methods. VQT shows high efficiency.**

| Type | TFLOPs | Device | Time | Speed Ratio |
|---|---|---|---|---|
| VMAF | - | CPU | 9.85s | 1.0× |
| AVQT | - | CPU | 4.61s | 2.1× |
| MDTSVFA | 231 | GPU | 7.07s | 1.4× |
| StarVQA | 197 | GPU | 0.57s | 17.3× |
| BVQA | 89 | GPU | 0.51s | 19.3× |
| STDAM | 106 | GPU | 2.12s | 4.6× |
| VQT | 154 | GPU | 0.50s | 19.7× |

obtained by combining them. We also verify individual modules on Kinetics-400, where a combination of modules improves 1.3% on Top-1 accuracy which proves the ability of generalization on general classification tasks. For better understanding, we plot the visualization of the learned temporal attention maps in Fig. 4. The visualization shows that different STA modules in the MPTN pay attention to different frames in a clip. It means that the proposed MPTN can capture different distortions simultaneously.

*Visualization.* To further validate the effectiveness of VQT, we visualize the frames with the highest response and analyze the corresponding quality scores for each individual frame. As shown in Fig. 5, 6 and 7, VQT is capable of effectively perceiving co-existing low-quality features in videos, such as interlace, motion blur, out-of-defocus and blocking artifacts. Furthermore, compared to averaging the predicted results for all frames, VQT-based temporal processing yields prediction results that are closer to the labeled MOS values.

*Efficiency comparison.* Computation Efficiency is compared with SoTA methods under 1080P/30FPS/30s videos, as shown in Tab. 5, including academic algorithms and industrial algorithms. Due to commercial confidentiality, we cannot obtain open-source model information (*e.g.*, , FLOPs) from Netflix/Apple. For a fair comparison, in GPU, VQT shows higher efficiency than image-based (MDTSVAF, STDAM) and video-based SoTA methods (StarVQA, BVQA). That 0.5s inference time can fulfill the real-time monitoring of the service-side quality variation. Further, speed-up can be investigated in future work (*e.g.*, , knowledge distillation and quantization).

## 4.6 Generalization in Video Classification

To further evaluate the generalization ability of VQT to other general semantic task, the performance on the video classification dataset of Kinetics-400[5] with SoTA video classification models is evaludated, including R(2+1)D [47], I3D [5], I3D+NL [52], ip-CSN-152 [46], SlowFast [15], TimeSformer [3] and Video Swin Transformer [35]. We follow the default training setting in MMAction2 for a fair comparison. As shown in Tab. 6, VQT achieves a Top-1 accuracy by 80.3%, outperforming CNN models [46, 47, 58] and very recent Transformer-based models, *e.g.*, TimeSformer[3] and Video Swin Transformer[35]. This proves the effectiveness and generalization ability of VQT in the general video classification task. We hope this VQT module can achieve more satisfactory performance when used in more general computer vision tasks.
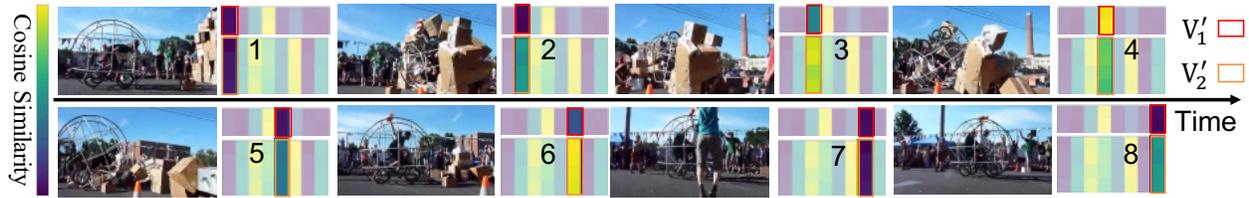
**Figure 4: Visualization results of the temporal attention maps generated by an MPTN consisting of two STA modules (No.5956265529 in the KoNViD-1k). By the comparison between $V'_1$ and $V'_2$, these two STA modules concentrate on different distortions. The STA with fewer frames pays much attention to compression distortions that can be detected by spatial information, and the STA with more frames focuses on camera movement that can be detected by temporal information.**
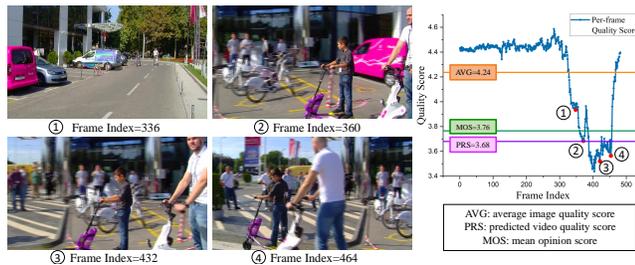


**Figure 5: Visualization results of the No.TelevisionClip-0604 video sampled from Youtube-UGC. The mixture of interlace and motion blur deteriorates the video quality. Compared to the average image quality score, VQT focuses more on the clips containing distortions, as shown by 4 sampled frames.**



**Figure 7: Visualization results of the No.0095_47 video sampled from ICME. The video quality is mainly decided by blocking artifacts, which appear in the fighting scenarios.**

**Table 6: Classification results on the K-400 validation set. The computational FLOPs and Top-k accuracy are reported.**

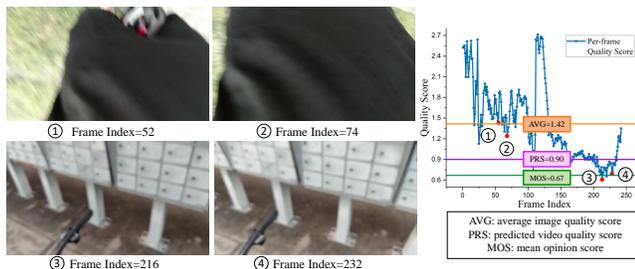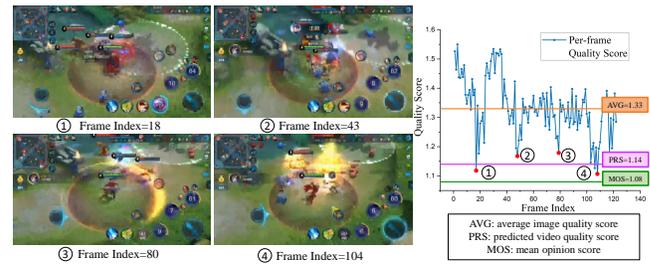| Method | Backbone | Top-1 | Top-5 | TFLOPs |
|---|---|---|---|---|
| R(2+1)D | ResNet34 | 72.0 | 90.0 | 75 |
| I3D | ResNet50 | 72.1 | 90.3 | 108 |
| I3D+NL | ResNet101 | 77.7 | 93.3 | 359 |
| ip-CSN-152 | Res152 | 77.8 | 92.8 | 109 |
| SlowFast | ResNet101 | 79.8 | 93.9 | 234 |
| TimeSformer | ViT-B | 78.0 | 93.7 | 197 |
| Video Swin | Swin-T | 78.8 | 93.6 | 88 |
| VQT | ViT-B | **80.3** | **94.5** | **154** |



**Figure 6: Visualization results of the No.B174 video sampled from LIVE-VQC. Our analysis indicates that the primary factors contributing to the degradation are motion blur (exhibited in the 1st and 2nd frames) and out-of-focus (exhibited in the 3rd and 4th frames).**

## 5 CONCLUSION AND FUTURE WORK

This paper proposed VQT to address two underestimated challenges faced by VQA. To tackle the first challenge that the perceptual quality of videos is largely determined by deteriorated keyframes, we propose the STA module, which performs sparse sampling by analyzing the correlation between frames, resulting in efficient computation of attention. To address the second challenge that various types of distortions co-exist in a video, we propose the MPTN capture co-existing distortions by stacking multiple STA

modules with different degrees of sparsity. Our proposed method is extensively evaluated on three widely-used NR-VQA datasets. Additionally, VQT outperforms widely-adopted industrial algorithms of VMAF and AVQT on four FR-VQA datasets. Extensive ablation studies and visual analysis further validate the effectiveness of each component of VQT. We also observe good generalization ability when transferring to the video classification task. We hope that VQT can serve as a new baseline for VQA tasks.

Regarding the selection of keyframes, the STA module currently relies on predefined hyperparameters. In the future, it would be possible to propose methods for adaptive keyframe selection, which can determine the number of keyframes based on prior information or employ a learning-based approach to identify keyframes that exhibit distortion-related features.

# REFERENCES

[1] Aishy Amer and Eric Dubois. 2005. Fast and reliable structure-oriented video noise estimation. *IEEE TCSVT* 15, 1 (2005), 113–118.

[2] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lucic, and Cordelia Schmid. 2021. ViViT: A Video Vision Transformer. In *ICCV*. IEEE, 6816–6826.

[3] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. 2021. Is Space-Time Attention All You Need for Video Understanding?. In *ICML*, Vol. 139. PMLR, 813–824.

[4] Sebastian Bosse, Dominique Maniry, Thomas Wiegand, and Wojciech Samek. 2016. A deep neural network for image quality assessment. In *ICIP*. IEEE, 3773–3777.

[5] João Carreira and Andrew Zisserman. 2017. Quo Vadis, Action Recognition? A New Model and the Kinetics Dataset. In *CVPR*. IEEE Computer Society, 4724–4733.

[6] Aaron Chadha and Yiannis Andreopoulos. 2021. Deep Perceptual Preprocessing for Video Coding. In *CVPR*. Computer Vision Foundation / IEEE, 14852–14861.

[7] Krzysztof Marcin Choromanski, Valerii Likhosherstov, David Dohan, Xingyou Song, Andreea Gane, Tamás Sarlós, Peter Hawkins, Jared Quincy Davis, Afroz Mohiuddin, Lukasz Kaiser, David Benjamin Belanger, Lucy J. Colwell, and Adrian Weller. 2021. Rethinking Attention with Performers. In *ICLR*.

[8] MMAction2 Contributors. 2020. OpenMMLab's Next Generation Video Understanding Toolbox and Benchmark. https://github.com/open-mmlab/mmaction2.

[9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. ImageNet: A large-scale hierarchical image database. In *CVPR*. IEEE Computer Society, 248–255.

[10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In *ICLR*.

[11] Joshua Peter Ebenezer, Zaixi Shang, Yongjun Wu, Hai Wei, and Alan C. Bovik. 2020. No-Reference Video Quality Assessment Using Space-Time Chips. In *MMSP*. IEEE, 1–6.

[12] Hossein Talebi Esfandarani and Peyman Milanfar. 2018. Learned perceptual image enhancement. In *ICCP*. 1–13.

[13] Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra Malik, and Christoph Feichtenhofer. 2021. Multiscale Vision Transformers. In *ICCV*. 6824–6835.

[14] Yuming Fang, Hanwei Zhu, Yan Zeng, Kede Ma, and Zhou Wang. 2020. Perceptual Quality Assessment of Smartphone Photography. In *CVPR*. Computer Vision Foundation / IEEE, 3674–3683.

[15] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. 2019. Slow-Fast Networks for Video Recognition. In *ICCV*. 6201–6210.

[16] Fei Gao, Yi Wang, Panpeng Li, Min Tan, Jun Yu, and Yani Zhu. 2017. DeepSim: Deep similarity for image quality assessment. *Neurocomputing* 257 (2017), 104–114.

[17] Deepti Ghadiyaram, Janice Pan, and Alan C. Bovik. 2019. A Subjective and Objective Study of Stalling Events in Mobile Streaming Videos. *IEEE Trans. Circuits Syst. Video Technol.* 29, 1 (2019), 183–197.

[18] Jie Gu, Gaofeng Meng, Cheng Da, Shiming Xiang, and Chunhong Pan. 2019. No-Reference Image Quality Assessment with Reinforcement Recursive List-Wise Ranking. In *AAAI*. AAAI Press, 8336–8343.

[19] Jinwei Gu, Ravi Ramamoorthi, Peter N. Belhumeur, and Shree K. Nayar. 2009. Removing image artifacts due to dirty camera lenses and thin occluders. *ACM Trans. Graph.* 28, 5 (2009), 144.

[20] Vlad Hosu, Franz Hahn, Mohsen Jenadeleh, Hanhe Lin, Hui Men, Tamás Szirányi, Shujun Li, and Dietmar Saupe. 2017. The Konstanz natural video database (KoNViD-1k). In *QoMEX*. IEEE, 1–6.

[21] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. 2016. Perceptual Losses for Real-Time Style Transfer and Super-Resolution. In *ECCV*. 694–711.

[22] William Johnson and Joram Lindenstrauss. 1984. Extensions of Lipschitz maps into a Hilbert space. *Contemp. Math.* 26 (1984).

[23] Le Kang, Peng Ye, Yi Li, and David S. Doermann. 2014. Convolutional Neural Networks for No-Reference Image Quality Assessment. In *CVPR*. IEEE Computer Society, 1733–1740.

[24] Will Kay, João Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, Mustafa Suleyman, and Andrew Zisserman. 2017. The Kinetics Human Action Video Dataset. *CoRR* abs/1705.06950 (2017).

[25] Junjie Ke, Qifei Wang, Yilin Wang, Peyman Milanfar, and Feng Yang. 2021. MUSIQ: Multi-scale Image Quality Transformer. In *ICCV*. 5128–5137.

[26] Jongyoo Kim and Sanghoon Lee. 2017. Deep Learning of Human Visual Sensitivity in Image Quality Assessment Framework. In *CVPR*. 1969–1977.

[27] Woojae Kim, Jongyoo Kim, Sewoong Ahn, Jinwoo Kim, and Sanghoon Lee. 2018. Deep Video Quality Assessor: From Spatio-Temporal Visual Sensitivity to a Convolutional Neural Aggregation Network. In *ECCV*, Vol. 11205. 224–241.

[28] Bowen Li, Weixia Zhang, Meng Tian, Guangtao Zhai, and Xianpei Wang. 2021. Blindly Assess Quality of In-the-Wild Videos via Quality-aware Pre-training and Motion Perception. *CoRR* abs/2108.08505 (2021).

[29] Dingquan Li, Tingting Jiang, and Ming Jiang. 2019. Quality Assessment of In-the-Wild Videos. In *ACM Multimedia*. ACM, 2351–2359.

[30] Dingquan Li, Tingting Jiang, and Ming Jiang. 2021. Unified Quality Assessment of in-the-Wild Videos with Mixed Datasets Training. *IJCV* 129 (2021), 1238–1257.

[31] Yuming Li, Lai-Man Po, Litong Feng, and Fang Yuan. 2016. No-reference image quality assessment with deep convolutional neural networks. In *DSP*. IEEE, 685–689.

[32] Zhuoran Li, Zhengfang Duanmu, Wentao Liu, and Zhou Wang. 2019. AVC, HEVC, VP9, AVS2 or AV1? - A Comparative Study of State-of-the-Art Video Encoders on 4K Videos. In *ICIAR (1) (Lecture Notes in Computer Science, Vol. 11662)*. Springer, 162–173.

[33] Min Liu, Guangtao Zhai, Zhenyu Zhang, Yuntao Sun, Ke Gu, and Xiaokang Yang. 2014. Blind image quality assessment for noise. In *BMSB*. 1–5.

[34] Yutao Liu, Ke Gu, Guangtao Zhai, Xianming Liu, Debin Zhao, and Wen Gao. 2017. Quality assessment for real out-of-focus blurred images. *JVCIR* (2017).

[35] Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. 2021. Video Swin Transformer. *CoRR* abs/2106.13230 (2021).

[36] Pina Marziliano, Frédéric Dufaux, Stefan Winkler, and Touradj Ebrahimi. 2002. A no-reference perceptual blur metric. In *ICIP*. 57–60.

[37] Xiongkuo Min, Ke Gu, Guangtao Zhai, Menghan Hu, and Xiaokang Yang. 2018. Saliency-induced reduced-reference quality index for natural scene and screen content images. *SP* 145 (2018), 127–136.

[38] Anish Mittal, Anush Krishna Moorthy, and Alan Conrad Bovik. 2012. No-Reference Image Quality Assessment in the Spatial Domain. *IEEE Trans. Image Process.* 21, 12 (2012), 4695–4708.

[39] Anish Mittal, Michele A. Saad, and Alan C. Bovik. 2016. A Completely Blind Video Integrity Oracle. *IEEE TIP* 25, 1 (2016), 289–300.

[40] Anish Mittal, Rajiv Soundararajan, and Alan C. Bovik. 2013. Making a "Completely Blind" Image Quality Analyzer. *IEEE SPL* 20, 3 (2013), 209–212.

[41] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Z. Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *NeurIPS*. 8024–8035.

[42] Michele A. Saad, Alan C. Bovik, and Christophe Charrier. 2014. Blind Prediction of Natural Video Quality. *IEEE Trans. Image Process.* 23, 3 (2014), 1352–1365.

[43] Zeina Sinno. 2019. Large-Scale Study of Perceptual Video Quality. *IEEE TIP* 28, 2 (2019), 612–627.

[44] Ahmed Telili, Sid Ahmed Fezza, Wassim Hamidouche, and Hanene FZ Meftah. 2022. 2BiVQA: Double Bi-LSTM based Video Quality Assessment of UGC Videos. *arXiv preprint arXiv:2208.14774* (2022).

[45] Bart Thomee, David A. Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li. 2016. YFCC100M: the new data in multimedia research. *Commun. ACM* 59, 2 (2016), 64–73.

[46] Du Tran, Heng Wang, Matt Feiszli, and Lorenzo Torresani. 2019. Video Classification With Channel-Separated Convolutional Networks. In *ICCV*.

[47] Du Tran, Heng Wang, Lorenzo Torresani, Jamie Ray, Yann LeCun, and Manohar Paluri. 2018. A Closer Look at Spatiotemporal Convolutions for Action Recognition. In *CVPR*. 6450–6459.

[48] Zhengzhong Tu, Yilin Wang, Neil Birkbeck, Balu Adsumilli, and Alan C. Bovik. 2021. UGC-VQA: Benchmarking Blind Video Quality Assessment for User Generated Content. *IEEE TIP* 30 (2021), 4449–4464.

[49] Zhengzhong Tu, Xiangxu Yu, Yilin Wang, Neil Birkbeck, Balu Adsumilli, and Alan C. Bovik. 2021. RAPIQUE: Rapid and Accurate Video Quality Prediction of User Generated Content. *CoRR* abs/2101.10955 (2021).

[50] Sinong Wang, Belinda Z. Li, Madian Khabsa, Han Fang, and Hao Ma. 2020. Linformer: Self-Attention with Linear Complexity. *CoRR* abs/2006.04768 (2020).

[51] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. 2021. Pyramid Vision Transformer: A Versatile Backbone for Dense Prediction without Convolutions. *CoRR* abs/2102.12122 (2021).

[52] Xiaolong Wang, Ross B. Girshick, Abhinav Gupta, and Kaiming He. 2018. Non-Local Neural Networks. In *CVPR*. 7794–7803.

[53] Yilin Wang and Balu Adsumilli. 2019. YouTube UGC Dataset for Video Compression Research. In *MMSP*. IEEE, 1–5.

[54] Yilin Wang, Junjie Ke, Hossein Talebi, Joong Gon Yim, Neil Birkbeck, Balu Adsumilli, Peyman Milanfar, and Feng Yang. 2021. Rich Features for Perceptual Quality Assessment of UGC Videos. In *CVPR*. 13435–13444.

[55] Yunbo Wang, Mingsheng Long, Jianmin Wang, and Philip S. Yu. 2017. Spatiotemporal Pyramid Network for Video Action Recognition. In *CVPR*. IEEE Computer Society, 2097–2106.

[56] Zhou Wang, Alan C. Bovik, and Brian L. Evans. 2000. Blind Measurement of Blocking Artifacts in Images. In *ICIP*. IEEE, 981–984.

[57] Haoning Wu, Chaofeng Chen, Liang Liao, Jingwen Hou, Wenxiu Sun, Qiong Yan, and Weisi Lin. 2022. DisCoVQA: Temporal Distortion-Content Transformers for Video Quality Assessment. *CoRR* abs/2206.09853 (2022).

Kun Yuan[†], Zishang Kong[†], Chuanchuan Zheng, Ming Sun, and Xing Wen

[58] Saining Xie, Chen Sun, Jonathan Huang, Zhuowen Tu, and Kevin Murphy. 2018. Rethinking Spatiotemporal Feature Learning: Speed-Accuracy Trade-offs in Video Classification. In *ECCV*. 318–335.

[59] Fengchuang Xing, Yuan-Gen Wang, Hanpin Wang, Leida Li, and Guopu Zhu. 2021. StarVQA: Space-Time Attention for Video Quality Assessment. *CoRR* abs/2108.09635 (2021).

[60] Jiahua Xu, Jing Li, Xingguang Zhou, Wei Zhou, Baichao Wang, and Zhibo Chen. 2021. Perceptual Quality Assessment of Internet Videos. In *ACM Multimedia*. ACM, 1248–1257.

[61] Ceyuan Yang, Yinghao Xu, Bo Dai, and Bolei Zhou. 2020. Video Representation Learning with Visual Tempo Consistency. *CoRR* abs/2006.15489 (2020).

[62] Peng Ye, Jayant Kumar, Le Kang, and David S. Doermann. 2012. Unsupervised feature learning framework for no-reference image quality assessment. In *CVPR*. IEEE Computer Society, 1098–1105.

[63] Zhenqiang Ying, Maniratnam Mandal, Deepti Ghadiyaram, and Alan C. Bovik. 2021. Patch-VQ: 'Patching Up' the Video Quality Problem. In *CVPR*. Computer Vision Foundation / IEEE, 14019–14029.

[64] Junyong You and Jari Korhonen. 2019. Deep Neural Networks for No-Reference Video Quality Assessment. In *ICIP*. IEEE, 2349–2353.

[65] Kun Yuan, Shaopeng Guo, Ziwei Liu, Aojun Zhou, Fengwei Yu, and Wei Wu. 2021. Incorporating Convolution Designs Into Visual Transformers. In *ICCV*. IEEE, 579–588.

[66] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. 2018. The Unreasonable Effectiveness of Deep Features as a Perceptual Metric. In *CVPR*. 586–595.

[67] Yu Zhang, Xinbo Gao, Lihuo He, Wen Lu, and Ran He. 2019. Blind Video Quality Assessment With Weakly Supervised Learning and Resampling Strategy. *IEEE TCSVT* 29, 8 (2019), 2244–2255.

[68] Kai Zhao, Kun Yuan, Ming Sun, Mading Li, and Xing Wen. 2023. Quality-Aware Pre-Trained Models for Blind Image Quality Assessment. In *CVPR*. IEEE Computer Society, 22302–22313.

[69] Kai Zhao, Kun Yuan, Ming Sun, and Xing Wen. 2023. Zoom-VQA: Patches, Frames and Clips Integration for Video Quality Assessment. In *CVPR Workshops*. IEEE Computer Society, 1302–1310.

[70] Luo-yu Zhou and Zheng-bing Zhang. 2014. No-reference image quality assessment based on noise, blurring and blocking effect. *Optik* 125, 19 (2014), 5677–5680.