# DisCoVQA: Temporal Distortion-Content Transformers for Video Quality Assessment

Haoning Wu, Chaofeng Chen, Liang Liao, *Member, IEEE*, Jingwen Hou, *Student Member, IEEE*, Wenxiu Sun, Qiong Yan, Weisi Lin, *Fellow, IEEE*

*Abstract*—The temporal relationships between frames and their influences on video quality assessment (VQA) are still under-studied in existing works. These relationships lead to two important types of effects for video quality. Firstly, some temporal variations (such as shaking, flicker, and abrupt scene transitions) are causing temporal distortions and lead to extra quality degradations, while other variations (e.g. those related to meaningful happenings) do not. Secondly, the human visual system often has different attention to frames with different contents, resulting in their different importance to the overall video quality. Based on prominent time-series modeling ability of transformers, we propose a novel and effective transformer-based VQA method to tackle these two issues. To better differentiate temporal variations and thus capture the temporal distortions, we design a transformer-based Spatial-Temporal Distortion Extraction (STDE) module. To tackle with temporal quality attention, we propose the encoder-decoder-like temporal content transformer (TCT). We also introduce the temporal sampling on features to reduce the input length for the TCT, so as to improve the learning effectiveness and efficiency of this module. Consisting of the STDE and the TCT, the proposed Temporal Distortion-Content Transformers for Video Quality Assessment (*DisCoVQA*) reaches state-of-the-art performance on several VQA benchmarks without any extra pre-training datasets and up to 10% better generalization ability than existing methods. We also conduct extensive ablation experiments to prove the effectiveness of each part in our proposed model, and provide visualizations to prove that the proposed modules achieve our intention on modeling these temporal issues. We will publish our codes and pretrained weights later.

*Index Terms*—Deep learning, video quality assessment, temporal modeling, transformers

## I. INTRODUCTION

WITH the rapid development of smartphones and portable cameras, more and more videos are created every day by a huge diversity of consumers, both professional or non-professional users. These videos are collected in-the-wild and uploaded to popular social media websites or apps such as YouTube and TikTok, and the number of them is still growing. Henceforth, it becomes increasingly important to build a robust and powerful objective VQA method for these natural videos without pristine references.

Many existing VQA efforts [20], [23], [32], [38], [40], [42] have demonstrated that directly adopting the image quality

H. Wu, C. Chen, L. Liao, J. Hou, and W. Lin are with the School of Computer Engineering, Nanyang Technological University, Singapore. (e-mail: haoning001@e.ntu.edu.sg; chaofeng.chen@ntu.edu.sg; liang.liao@ntu.edu.sg; jingwen003@e.ntu.edu.sg; wslin@ntu.edu.sg;)

W. Sun and Q. Yan are with the Sensetime Research, Hong Kong, China. (e-mail: irene.wenxiu.sun@gmail.com; sophie.yanqiong@gmail.com)

Corresponding author: Weisi Lin.

assessment (IQA) models to measure video quality with frame-wise IQA procedures is not effective for VQA, as they neglect the temporal relationships between frames. Specifically, TLVQM [20] have demonstrated that the variations between frames, such as shaking, flicker and abrupt scene transitions, are causing additional quality degradations. Such effects have also been observed by many deep VQA approaches such as CoINVQ [40] and PVQ [42]. We categorize all these degradations caused by variations between different frames as ***temporal distortions***, the first type of temporal issues in VQA. While these distortions are usually noticed, many approaches at present still apply handcrafted kernels to model them. Despite temporal distortions, many existing approaches have also noticed that different frames should have different importance to the final video quality. VSFA [23] and some other recent approaches [22], [24] have concluded this effect as "temporal memory effect" that early frames are prone to be forgotten, yet recent studies on the human visual system (HVS) such as [10] imply that importance of frames should be highly related to their contents, so we categorize this effect as ***content-related temporal quality attention*** instead. These two temporal issues are often mentioned in prior arts but less systematically discussed or modeled. Recently, transformer architectures are proved to have better ability on time-series modeling in several tasks, including language modeling [39] and video recognition [1], [8], [25], which allows our method to revisit these two issues discussed above with powerful transformers. In the following two paragraphs, we will discuss them in detail and provide our corresponding approaches based on transformers.

First of all, for temporal distortions, it is straightforward to notice that they come from temporal variations between frames, as static videos without variations do not have temporal distortions. As illustrated in Fig. 1(a), some temporal variations, such as the shaking of the whole picture in *video 1*, the rapid brightness changes among frames in *video 2* and the abrupt transition of the scene in *video 3*, do not have clear semantics or reason and lead to temporal distortions. On the other hand, not all temporal variations lead to temporal distortions. To be specific, some other temporal variations come from human actions or other meaningful activities/events in videos, such as running players and correspondingly moving scene during a football game (*video 4*), and therefore do not lead to distortions in general. The aforementioned two scenarios suggest that understanding the semantics of actions (or other activities) helps to better capture temporal distortions. Therefore, we adopt the Video Swin Transformer Tiny
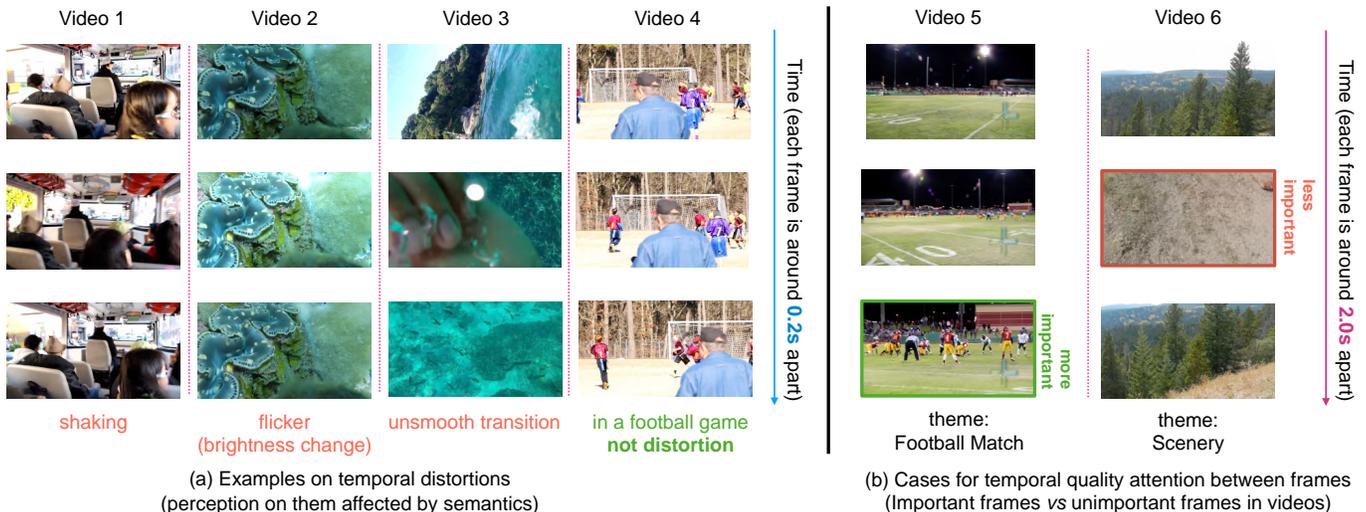
Fig. 1: Examples for naturally collected videos in LSVQ dataset [42] show complicated temporal relationships and their effects to VQA. These effects have shown that better temporal modeling are in need for VQA task.

(*abbr.* as Swin-T) with better action recognition ability as the backbone for our Spatial-Temporal Distortion Extraction (STDE). Such semantic information cannot be captured by existing classical handcraft models [20], [32], [38] and even 3D-CNNs [4], [14], [36] as introduced in [42], [44], which are less effective than transformers on several action-related tasks as compared in [2], [8], [25]. In addition, we apply the temporal differences to further enhance the sensitivity for temporal variations, especially for those variations independent to actions (*e.g.* scene transitions) and are less captured by the backbone network. With Swin-T and temporal differences, we better extract features that are sensitive to distortion-related temporal variations in STDE.

Besides temporal distortions, the temporal quality attention mechanism also affects the overall quality evaluation of the video. Similar to spatial quality attention that is related to contents of different regions [18], [34], the temporal attention should also be related to contents of frames. For example, as shown in Fig. 1(b), *video 5*, in the replay video of a football match, the attention of the HVS is more attracted by the frames that contain the zoomed-in players (in the green box), which are closer to what the whole video is about, *i.e.* the video theme. On the contrary, some frames less related to the video theme might be less important in deciding the final video quality. For example, the poor quality of the intermediate frames (in the red box) shot to the ground does not affect the overall video quality of *video 6* rated as good (as scored by the human in [42]), as the HVS pays less attention on such frames but more to those frames about *Scenery* (theme of *video 6*). Both examples suggest that the relevance of frame contents to the overall video theme affects the importance of frames in deciding the final video quality. As transformers are well recognized as good at learning correlations among a sequence [7], [39], we propose the temporal content transformer (TCT) to learn the correlations of frame contents and model this temporal quality attention. The TCT has an encoder-decoder-like transformer structure

that better extracts the temporal quality attention from frame contents. Some existing approaches [5], [23], [24], [44] also introduce RNNs such as GRU [6] or LSTM [16] to model the temporal quality attention. However, RNN-based models are usually weak in modeling long-range correlation, so they are less effective in extracting the correlations between frames across the whole video and their relevance to the overall theme, which is better modeled by the proposed TCT.

As analyzed above, we introduced two transformer-based modules, the STDE and the TCT, to improve temporal modeling in VQA. To better apply transformers in VQA, we present two vital designs to improve the effectiveness of the two transformer-based modules. For the STDE, as the last-layer output of the transformer backbone might be less sensitive to low-level features and temporal variations, we introduce multi-level feature extraction to capture both low-level and high-level features and thus better spot temporal distortions. For the TCT, the transformer-based architecture is hard to be learnt effectively when datasets are small or input sequences are too long. Also, the long inputs lead to quadratically-increased computational costs. To improve the training effectiveness and efficiency, we propose the temporal sampling on features (TSF) to reduce the input length for the TCT during training. By cutting video features into segments and randomly sampling one feature frame from each segment, the TSF significantly improves the performance of training the TCT on small or long-duration video datasets. With the STDE, the TCT, and the vital designs discussed above, we build an efficient and effective transformer-based VQA model, as illustrated in Fig. 2.

In summary, we discuss the temporal relationships in VQA and propose the **T**emporal **Dis**tortion-**Co**ntent Transformers for **V**ideo **Q**uality **A**ssessment (**DisCoVQA**) that reaches state-of-the-art performance and good generalization ability on most natural VQA datasets, with the following contributions.

- We propose the Spatial-Temporal Distortion Extraction (STDE) to model the temporal distortions with consid-
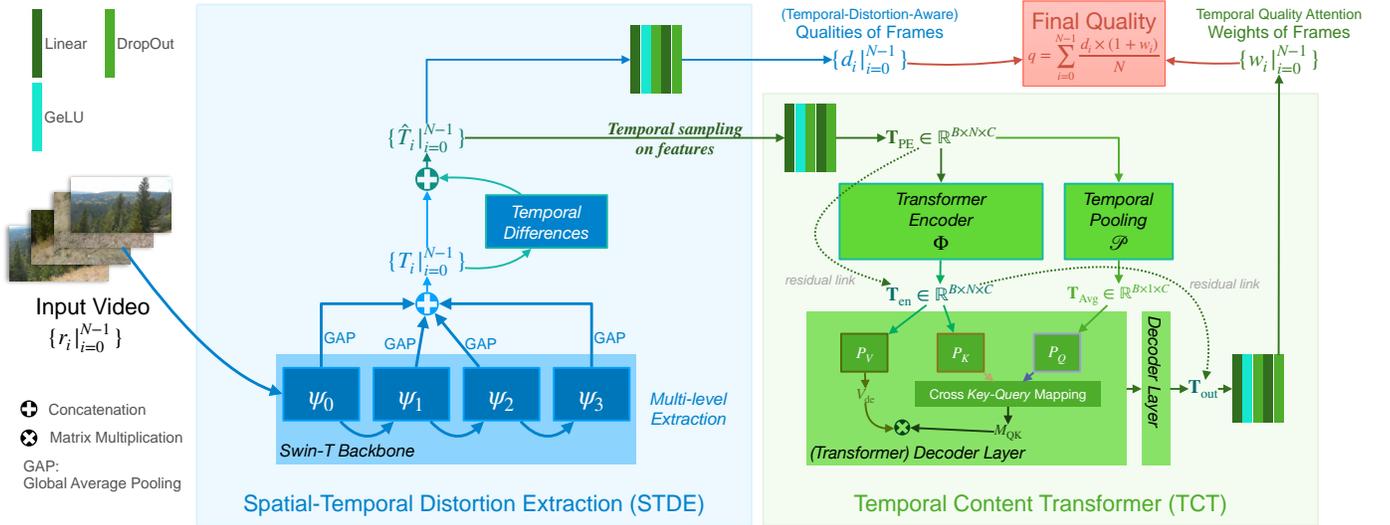
Fig. 2: The structure of proposed Temporal Distortion-Content Transformers for Video Quality Assessment (**DisCoVQA**). It contains the Spatial-Temporal Distortion Extraction (STDE, Sec. III-B) to better extract temporal distortions (such as *shaking, flicker, unsmooth transitions*), and the Temporal Content Transformer (TCT, Sec. III-C) to learn the content-related temporal quality attention between frames.

eration of the impact of semantics of actions in videos. It includes a Swin-T backbone and computes temporal differences to better model distortion-related temporal variations during the feature extraction phase in VQA.

- We propose the Temporal Content Transformer (TCT) to learn the temporal quality attention, considering the correlation of frames with the overall video theme. It includes a transformer encoder and a transformer decoder to learn quality attention weights for frames.
- We propose two important designs to adapt transformers to VQA: multi-level extraction to enhance low-level sensitivity for the STDE; and temporal sampling on features (TSF) to reduce the input length and improve the training effectiveness for the TCT. They both improve the final performance of the proposed DisCoVQA.

In the rest of the paper, we present the related work in Section II and the proposed temporal distortion-content transformer (DisCoVQA) in Section III. In section IV, we elaborate experimental settings and present extensive experimental results. In section V, we conduct the ablation studies, qualitative results and reliability analysis to explain its effectiveness. Finally, conclusions are drawn in Section VI.

## II. RELATED WORKS

In this section, we briefly review related works in each of the three sub-fields: classical video quality assessment (VQA) methods, prior deep VQA methods, and transformers on temporal modeling.

### A. Classical VQA methods

Classical VQA methods design on spatial [11], [27], [30] or spatial-temporal [20], [28], [32], [38] handcrafted features and regress on these features for quality modeling. Among

them, TLVQM [20] introduces a combination of spatial high-complexity and temporal low-complexity handcraft features and has reached state-of-the-art performance on LIVE-VQC [33]. VIDEVAL [38] ensembles different handcraft features and achieves better performance on KoNViD-1k [17]. Several most recent models, such as RAPIQUE [37] and CNN-TLVQM [21], use handcraft features for temporal modeling and deep neural networks for spatial modeling and reach far better results than pure spatial modeling approaches. Classical VQA methods have proved that the temporal relationships are non-negligible when evaluating video quality.

### B. Deep VQA Methods

During recent years, deep video quality assessment (VQA) methods have become predominant in VQA. Instead of extracting handcraft features, deep VQA methods [5], [13], [23] extract rich semantic features with CNN and regress the extracted features to predict video quality. For example, MLSP-FF [13] runs a linear regression on features extracted from video frames with Inception-ResNet-v2 [35]. These deep VQA methods also benefits from better temporal modeling. For instance, to model temporal distortions more effectively, several works [40], [42], [44] also introduce 3D-CNN instead of 2D-CNN to extract features. Some other deep VQA methods also use recurrent neural network (RNN) layers as temporal modeling for VQA, aiming to capture the temporal quality attention. For example, VSFA [23] uses ResNet-50 backbone [15] and GRU [6] as temporal regression, while some others [40], [44] also use LTSM [16], another type of RNN instead of GRU. LSCT-PHIQ [43] introduces naive transformer encoders for temporal attention modeling. A most recent model, *R+S* [22], combines CNN-based temporal distortion modeling and RNN-based temporal quality attention modeling together and reaches good performance. The existing

practice of deep VQA methods has further demonstrated the importance of temporal distortion representation and temporal quality attention modeling, which can significantly improve the prediction accuracy for VQA.

### C. Transformers on Temporal Modeling

Transformers proposed by [39] have been proved to have a powerful self-attention mechanism on time-series modeling. Many representative works such as BERT [7] and GPT-3 [3] have shown transformers can be far better on extracting attention from long sequences than RNNs ( [6], [16]) on language tasks. Transformers have also boosted video-related vision tasks. Many transformer-based video backbones, *e.g.* ViViT [1], MViT [8] and Video Swin Transformers [25], have achieved even better performance on video recognition with pure videos than conventional CNN-based methods combined with optical flows and much better than pure 3D-CNN methods. As optical flows are originally designed to explicitly compute temporal variations, the success of video transformer backbones suggests that they are already strong on modeling temporal variations. All these existing practices have suggested transformers are stronger structures for different types of temporal modeling.

## III. PROPOSED METHOD

### A. Framework Overview

The proposed DisCoVQA adopts a temporal hierarchical pipeline, as shown in Fig. 2. In DisCoVQA, we first input a video into the STDE (Sec. III-B). The STDE extracts feature tokens that are sensitive to temporal distortions (shaking, flicker, unsmooth transitions). Then the tokens are sampled by the temporal sampling on features (TSF) and passed into the TCT (Sec. III-C) to further model their content-related temporal quality attention. These attention weights are multiplied with the temporal-distortion-aware qualities of frames regressed by the STDE to obtain the overall quality prediction (Sec. III-D). Each of these components will be explained below.

### B. Spatial-Temporal Distortion Extraction

The Spatial-Temporal Distortion Extraction (STDE) is designed to better capture inter-frame temporal distortions during the feature extraction stage of VQA. It first extracts multi-level features from the Swin-T backbone, and compute temporal differences on extracted features, and then regress these features into temporal-distortion-aware qualities of frames. We explain each part of the STDE as follows.

*1) Multi-level Feature Extraction on Swin-T Backbone:* In STDE, we first extract features with a transformer-based backbone. For evaluation fairness, we choose the Video Swin Transformer Tiny (*abbreviated as* Swin-T) instead of its heavier versions which has similar parameters with I3D-ResNet-50 [4] (the most common ResNet-50 variant for videos) as our backbone. The Swin-T consists of four hierarchical Swin Transformer blocks ($\Psi_l, l = 0, 1, 2, 3$), where each block has several alternate 3D window multihead self-attention

(3DW-MSA) layers and feed-forward layers. A video clip $\mathcal{R} = \{r_i|_{i=0}^{N-1}\}$ is passed into Swin-T, where $i$ is the index for the $i$-th frame. Then the feature set $\mathcal{M}^l = \{M_i^l|_{i=0}^{N-1}\}$ for each frame of the video clip $\mathcal{R}$ is obtained from $l+1$ cascading blocks of Swin-T:

$$\mathcal{M}^l = \Psi_l(...(\Psi_0(\mathcal{R}))), \quad l = 0, 1, 2, 3, \quad (1)$$

Though adopting $l = 3$ is enough for recognition tasks, low-level quality-related information may not be sufficiently preserved in the resulting features. Therefore, we introduce the multi-level feature extraction on the Swin-T backbone to enhance sensitivity of features on low-level distortion-related information. During multi-level feature extraction, the features from different levels of Swin-T $\mathcal{M}^l$ are spatially pooled by the global average pooling layer (GAP), and concatenated into the primary tokens:

$$\mathcal{T} = \bigoplus_{l=0}^{3} \text{GAP}(\mathcal{M}^l), \quad (2)$$

where $\mathcal{T} = \{T_i|_{i=0}^{N-1}\}$ are the primary tokens for the video, and $\bigoplus$ denotes the concatenation operation on features from different levels of blocks. The multi-level feature extraction assures that both low-level features that contains distortion-related information and high-level semantics which is useful for Temporal Content Transformer in Sec. III-C are both captured by $\mathcal{T}$.

*2) Temporal Differences:* As discussed above, all these temporal distortions come from temporal variations. Though we introduce action recognition backbones to sense these variations, some other distortion-related variations such as unsmooth scene transitions are less related to video actions and are less captured by the backbone. Inspired by [19], [20], we extract the temporal differences between features of adjacent frames to further catch the temporal variations and then concatenate them with primary tokens $\mathcal{T}$ to get the STDE tokens $\hat{\mathcal{T}} = \{\hat{T}_i|_{i=0}^{N-1}\}$, defined as follows:

$$\hat{T}_i = T_i \oplus (T_i - T_{i+1}), \qquad 0 \le i < N - 1 \quad (3)$$

$$\hat{T}_i = T_i \oplus \mathbf{0}, \qquad i = N - 1 \quad (4)$$

where $\oplus$ denotes that two feature vectors are concatenated one after another in the channel dimension.

*3) Temporal-Distortion-Aware Qualities for Frames:* After the token extraction, we build a direct path with a two-layer multi-layer perceptron (MLP) to reduce STDE tokens $\{\hat{T}_i|_{i=0}^{N-1}\}$ into temporal-distortion-aware qualities for frames ($\{d_i|_{i=0}^{N-1}\}$). Denoting the two linear layers of the MLP as $l_1$ and $l_2$, and the GELU activation function as $A_{\text{GELU}}$, the $d_i$ is generated through the following equation:

$$d_i = l_2(A_{\text{GELU}}(l_1(T_i))) \quad (5)$$

The $\{\hat{T}_i|_{i=0}^{N-1}\}$ will also be fed into the TCT to get the temporal quality attention weights and combine with $\{d_i|_{i=0}^{N-1}\}$ to get the final video quality (as in Eq. (14). The structure and pipeline of the TCT is explained in the next section.

We visualize the temporal-distortion-aware frame qualities $\{d_i|_{i=0}^{N-1}\}$ in Sec. V-C.
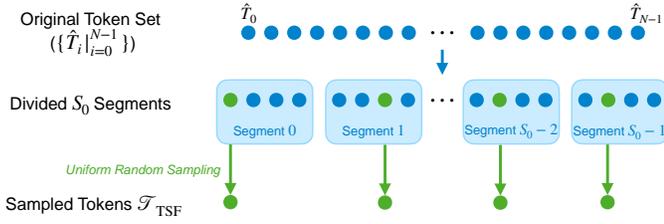
Fig. 3: The paradigm of proposed temporal sampling on features (TSF), where the features are divided into $S_0$ uniform segments. We sample multiple times for TSF during inference to improve its stability.

## C. Temporal Content Transformer

We design the temporal content transformer (TCT) to model the content-aware temporal quality attentions of frames with transformers. The TCT first conducts temporal sampling on features (TSF) from the STDE-extracted tokens $\{\hat{T}_i|_{i=0}^{N-1}\}$ and then reduces the channels of these features. These processed features are passed into a transformer encoder for correlation modeling between different frames, and then a transformer decoder to further learn the correlation of frame content to the overall theme. Finally, an MLP regresses the outputs of the transformers to get the temporal quality attention weights of frames. We discuss each part of the TCT as follows.

*1) Temporal Sampling on Features (TSF):* As transformers need matrix multiplications which is with $O(N^2)$ complexity to sequence length $N$, we use proposed temporally-sampled features instead of the original full features to reduce the input sequence length for the TCT. As illustrated in Fig. 3, we first divide the video into $S_0$ uniform-length segments regardless of the its original length and randomly sample one token from each segment. Given the input $\mathcal{T} = \{\hat{T}_i|_{i=0}^{N-1}\}$ (tokens extracted from STDE), the token set in temporal sampling on features $\mathcal{T}_{\text{TSF}}$ is expressed as:

$$\mathcal{T}_{\text{TSF}} = \{\hat{T}_{\mathcal{U}_{\text{int}}(\frac{j \times N}{S_0}, \frac{(j+1) \times N}{S_0})}|_{j=0}^{S_0-1}\}, \qquad (6)$$

where $\mathcal{U}_{\text{int}}(a, b)$ (where $a < b$) uniformly samples an integer within the range bounded by $a$ and $b$.

Contrary to the common concern that such sampling will reduce the performance, we notice that for VQA this practice improves the performance in some cases. This is because unlike natural language that is highly abstract, videos are usually continuous between frames and adjacent frames are usually with similar contents. Therefore, the TSF won't lose too much information about the video contents. On the contrary, compared with the vanilla transformer that uses all feature tokens as input, the TSF significantly reduces the token length and thus reduce the training difficulty of the transformer model. Benefited from its randomness, it also increases the total number of training pairs. This is especially helpful and improves the training effectiveness on small datasets such as CVD2014 [29] and LIVE-Qualcomm [12] where the TCT tend to overfit, or YouTube-UGC [41] dataset where the inputs are too long. (In Tab. XII) The TSF also significantly boosts the training speed and slightly improves the performance on other datasets (In Tab. XIII). Compared with temporal pooling

on inputs, it also keeps the original tokens free from any pooling kernels so that the quality information of frames is not corrupted during the processing. [1]

*2) Residual Transformer Encoder:* We introduce a transformer encoder on the TSF-sampled features to learn the content correlations of frames across the whole video. The core module in the transformer encoder is the self-attention module [39]. Generally, given a sequence of tokens $\mathbf{T} \in \mathbb{R}^{N \times C}$ [2] as input, the self-attention module will first project $\mathbf{T}$ into *key, query, value* matrices ($K, Q, V \in \mathbb{R}^{N \times C}$) with matrix multiplications by weights $P_K$, $P_Q$ and $P_V$, as follows:

$$K = \mathbf{T}P_K, Q = \mathbf{T}P_Q, V = \mathbf{T}P_V \qquad (7)$$

Then, $Q$ will multiply with the transposed $K^T \in \mathbb{R}^{C \times N}$ to get the $M \in \mathbb{R}^{N \times N}$ as follows:

$$M = \text{Softmax}(\frac{QK^T}{\sqrt{C}}) \qquad (8)$$

and $M^{i,j}$ is the attention value between element $i$ and $j$ in the sequence, reflecting the correlation between them. The computation of the attention value is agnostic to the distance between them and thus especially suitable to model global dependencies for temporal quality attention modeling in VQA.

The proposed transformer encoder $\Phi$ contains four sequential layers. Each layer includes a self-attention module as discussed in Eq. (7) and Eq. (8) and several feed-forward MLP layers, following the original structure as proposed in [39]. We also add a long-range residual link across the transformer encoder to enhance its learning effectiveness, and the whole residual transformer encoder that gets encoded tokens $\mathbf{T}_{\text{en}} \in \mathbb{R}^{N \times C}$ from the pre-encoding tokens $\mathbf{T}_{\text{pe}} \in \mathbb{R}^{N \times C}$ is expressed as follows:

$$\mathbf{T}_{\text{en}} = \Phi(\mathbf{T}_{\text{pe}}) + \mathbf{T}_{\text{pe}} \qquad (9)$$

*3) Transformer Decoder:* We carefully design a two-layer transformer decoder to detect specific frame contents that catch most human attention in the video, by referring the correlation of a specific frame's content with the average content. As this attention mechanism is related to the correlation of frame content to the overall topic (or theme) of the video, we design to explicitly capture these contents via the cross *key-query* mapping between the encoded token $\mathbf{T}_{\text{en}}$ and the average content token $\mathbf{T}_{\text{avg}}$. This cross *key-query* mapping is slight different from Eq. (7) and Eq. (8) and formulated as follows.

As the *Decoder Layer* in Fig. 2 shows, we perform *temporal pooling* ($\mathcal{P}_t$) on pre-encoding token to get the average token

$$\mathbf{T}_{\text{avg}} = \mathcal{P}_t(\mathbf{T}_{\text{pe}}) \in \mathbb{R}^{1 \times C} \qquad (10)$$

as the representation of the overall content. $\mathbf{T}_{\text{avg}}$ is projected into the *query* matrix. The *key* and *value* matrices are both projected from the encoded tokens $\mathbf{T}_{\text{en}} \in \mathbb{R}^{N \times C}$. Then this

---

[1] During inference, we randomly sample the tokens for $S_m$ multiple times and get the average result of different samples together to improve the prediction stability. Details and analyses can be found in Sec. V-D.

[2] where $C$ is the channel number and $N$ is the number of input tokens.

cross *key-query* mapping gives the $M_{QK}$ matrix computed as Eq. (12).

$$K_{de} = \mathbf{T}_{en}P_K, Q_{de} = \mathbf{T}_{avg}P_Q, V_{de} = \mathbf{T}_{en}P_V \quad (11)$$

$$M_{QK} = \text{Softmax}(\frac{Q_{de}K_{de}^T}{\sqrt{C}}) \quad (12)$$

The attention weights $M_{QK}$ in the last layer are directly multiplied with the *value* matrix projected from $\mathbf{T}_{en}$ and added with the $\mathbf{T}_{en}$ with a residual link (similar as in Eq. (9)) to get the output token $\mathbf{T}_{out}$ of the transformer decoder, which is further reduced to quality attention weights of frames $\{w_i|_{i=0}^{N-1}\}$ with the following scheme (denote the linear layers as $l_3, l_4$).

$$w = l_4(\text{A}_{\text{GELU}}(l_3(T_{en} + M_{QK}V_{de}))) \quad (13)$$

We visualize the learnt $\{w_i|_{i=0}^{N-1}\}$ and $M_{QK}$ in Sec. V-C.

### D. Final Video Quality Prediction

The attention weights $w_i$ are multiplied with the disortion-based qualities $d_i$ frame-by-frame and get the final video quality $q$ as follows:

$$q = \sum_{i=0}^{N-1} \frac{d_i \times (1 + w_i)}{N} \quad (14)$$

Following VQEG's suggestions and practices of several existing VQA works [22], [23], we remap the $q$ with consideration of subjective mean opinion scores (MOS) $s$ as follows:

$$\hat{q} = \frac{\max(s) - \min(s)}{1 + e^{\frac{q-\bar{q}}{\sigma(q)}}} + \min(s) \quad (15)$$

so that the predicted quality scores can be converted to the same range with the corresponding subjective scores. After remapping, we use the MAE loss between $\hat{q}$ and $s$

$$L = \|\hat{q} - s\|_1 \quad (16)$$

as our training loss function.

## IV. BENCHMARK EXPERIMENTS

In this section, we benchmark the performance of Dis-CoVQA on several natural VQA datasets. The proposed Dis-CoVQA shows the best performance in most individual dataset evaluation (Sec. IV-A4), and far better generalization ability in cross-dataset experiments (Sec. IV-C). DisCoVQA trained with large VQA datasets LSVQ and KoNViD-150k not only outperforms previous approaches but also shows competitive performance on small benchmark datasets without any fine-tuning process (Sec. IV-D). DisCoVQA is also with similar speed to non-transformer baselines (Sec. IV-E).

TABLE I: Sizes and characteristics of common VQA datasets.

| Dataset | Type | Distortion Type | Size |
|---|---|---|---|
| CVD2014 [29] | Normal Scale VQA dataset | Simulated Natural | 234 |
| LIVE-Qualcomm [12] | Normal Scale VQA dataset | Simlued Natural | 208 |
| KoNViD-1k [17] | Normal Scale VQA dataset | In-the-wild | 1,200 |
| LIVE-VQC [33] | Normal Scale VQA dataset | In-the-wild | 585 |
| LSVQ [42] | Large Scale VQA dataset | In-the-wild | 39,076 |
| KoNViD-150k [13] | Large Scale VQA dataset | In-the-wild | 150,000 |

### A. Experimental Settings

*1) Implementation Details:* We use one NVIDIA Tesla V100 GPU and Pytorch [31] for training. We set batch size $B = 16$, learning rate $lr = 0.001$ and use the AdamW [26] optimizer with $0.001$ weight decay rate during training. For multi-level feature extraction, we pass original-resolution videos without rescaling them. Following existing works [23], [24], [32], [38], we use SROCC (Spearman Rank-order Correlation Coefficients), PLCC (Pearson Linear Correlation Coefficients), KROCC (Kendall Rank-order Correlation Coefficients) as our evaluation metrics. Among them, SROCC & KROCC reflect the rank correlation of two sequences, where PLCC reflects the linear correlation, and higher correlations means better performances.

*2) Datasets:* We use six in-the-wild natural VQA datasets, where the majority of videos are directly photographed, instead of generated by algorithms, to benchmark the performance of our model (as listed in Tab. I). We select the first four, CVD2014, LIVE-Qualcomm, LIVE-VQC and KoNViD-1k with normal scale to benchmark the effectiveness of our model. These four datasets are further divided into two groups: the datasets with **simulated distortions** that are common in natural photography, CVD2014 and LIVE-Qualcomm; and the datasets **collected from real-world** natural videos, LIVE-VQC and KoNViD-1k. We report the *weighted average* performance of each group with respect to their dataset sizes to avoid the random biases during single dataset collections. For the rest two large scale VQA methods proposed recently, we also compare our model with existing approaches on them and evaluate whether directly training on these models can generalize well on other datasets.

*3) Baseline Methods:* We choose several most recent state-of-the-art methods (and label their time of publication) as comparison. Specifically, we compare with methods that represent different existing temporal modeling strategies in VQA, including VSFA [23], which applied a ResNet-50 2D-CNN backbone and an RNN for temporal modeling (and GST-VQA [5] which is based on VSFA and improves the training strategy for VQA); we also compare with CNN-TLVQM [21], which carefully designed handcrafted features for temporal modeling. We also notice a newly proposed method, MLSP-FF [13] with a heavy CNN backbone and only used naive average pooling for temporal modeling. A very recent approach, *R+S* applies SlowFast [9], a two-branch 3D-CNN network and a GRU [6] temporal regression module for temporal modeling, and ensembles it with another spatial branch for VQA.

*4) Evaluation Settings:* We conduct experiments on CVD2014 [29], LIVE-Qualcomm [12], LIVE-VQC [33] and KoNViD-1k [17], four individual VQA benchmark datasets

TABLE II: Result on CVD2014 [29], LIVE-Qualcomm [12], KoNViD-1k [17] and LIVE-VQC [33] datasets with standard 6:2:2 split setting. Some methods with original results not in this setting are reproduced by us. *Weighted Average* shows the weighted averaged result with respect to dataset sizes and the red and blue results represent the first and second best in chart. Methods with extra-dataset pre-training are labeled with (*EX*), and results that are neither reported nor reproduced are labeled as NA.

| Simulated Natural Datasets | CVD2014 (*234*, [29]) | | | LIVE-Qualcomm (*208*, [12]) | | | Weighted Average (442) | | |
|---|---|---|---|---|---|---|---|---|---|
| Method | SROCC(std) | PLCC(std) | KROCC(std) | SROCC(std) | PLCC(std) | KROCC(std) | SROCC | PLCC | KROCC |
| V-BLIINDS [32](2012,TIP) | 0.746(0.056) | 0.753(0.053) | 0.562(0.057) | 0.710(0.031) | 0.704(0.030) | 0.519(0.026) | 0.705 | 0.709 | 0.515 |
| TLVQM [20](2019,TIP) | 0.830(0.040) | 0.850(0.040) | NA | 0.780(0.070) | 0.810(0.060) | NA | 0.806 | 0.831 | NA |
| VSFA [23](2019,MM) | 0.870(0.037) | 0.868(0.032) | 0.695(0.047) | 0.773(0.061) | 0.795(0.055) | 0.587(0.062) | 0.824 | 0.834 | 0.644 |
| 3D-CNN+LSTM [44](2019,ICIP) | NA | NA | NA | 0.687 | 0.792 | NA | NA | NA | NA |
| MLSP-FF [13](2021,Access) | 0.770(0.060) | NA | NA | 0.710(0.080) | NA | NA | 0.742 | NA | NA |
| CNN-TLVQM [21](2020,MM) | 0.863(0.037) | 0.880(0.025) | 0.677(0.035) | 0.810(0.045) | 0.833(0.029) | 0.629(0.042) | 0.838 | 0.858 | 0.654 |
| GST-VQA [5](2021,TCSVT) | 0.831(0.052) | 0.844(0.062) | 0.657(0.037) | 0.801(0.053) | 0.825(0.043) | 0.620(0.052) | 0.817 | 0.835 | 0.639 |
| *R+S*(*EX*) [22](2022,TCSVT) | 0.860(0.037) | 0.877(0.034) | 0.687(0.048) | 0.814(0.045) | 0.819(0.054) | 0.639(0.057) | 0.838 | 0.849 | 0.664 |
| **DisCoVQA (Ours)** | 0.897(0.025) | 0.893(0.025) | 0.726(0.033) | 0.823(0.033) | 0.825(0.030) | 0.645(0.033) | 0.862 | 0.862 | 0.688 |
| Real-world Natural Datasets | LIVE-VQC (*585*, [33]) | | | KoNViD-1k (*1200*, [17]) | | | Weighted Average (1785) | | |
| Method | SROCC(std) | PLCC(std) | KROCC(std) | SROCC(std) | PLCC(std) | KROCC(std) | SROCC | PLCC | KROCC |
| V-BLIINDS [32](2012,TIP) | 0.694(0.050) | 0.718(0.050) | 0.508(0.042) | 0.710(0.031) | 0.704(0.030) | 0.519(0.026) | 0.705 | 0.709 | 0.515 |
| TLVQM [20](2019,TIP) | 0.799(0.036) | 0.803(0.036) | 0.608(0.037) | 0.773(0.024) | 0.768(0.023) | 0.577(0.022) | 0.782 | 0.779 | 0.587 |
| VIDEVAL [38](2021,TIP) | 0.752(0.039) | 0.751(0.042) | 0.564(0.036) | 0.783(0.021) | 0.780(0.021) | 0.585(0.021) | 0.773 | 0.770 | 0.578 |
| VSFA [23](2019)(2019,MM) | 0.773(0.027) | 0.795(0.026) | 0.581(0.031) | 0.773(0.019) | 0.775(0.019) | 0.578(0.019) | 0.773 | 0.782 | 0.579 |
| CNN-TLVQM [21](2020,MM) | 0.814(0.027) | 0.821(0.025) | 0.622(0.033) | 0.816(0.018) | 0.818(0.019) | 0.626(0.018) | 0.815 | 0.819 | 0.625 |
| RAPIQUE [37](2021,OJSP) | 0.755 | 0.786 | NA | 0.803 | 0.817 | NA | NA | NA | NA |
| MLSP-FF [13](2021,Access) | 0.720(0.060) | NA | NA | 0.820(0.020) | NA | NA | 0.787 | NA | NA |
| CoINVQ(*EX*) [40](2021,CVPR) | NA | NA | NA | 0.767 | 0.764 | NA | NA | NA | NA |
| PVQ(*EX*) [42](2021,CVPR) | 0.803(0.029) | 0.811(0.028) | 0.616(0.031) | 0.785(0.021) | 0.774(0.028) | 0.576(0.020) | 0.791 | 0.786 | 0.589 |
| LSCT-PHIQ(*EX*) [43](2021,MM) | 0.796(0.025) | 0.782(0.024) | 0.589(0.023) | 0.833(0.027) | 0.834(0.024) | 0.638(0.019) | 0.821 | 0.817 | 0.625 |
| GST-VQA [5](2021,TCSVT) | 0.788(0.032) | 0.796(0.028) | 0.604(0.037) | 0.814(0.026) | 0.825(0.043) | 0.621(0.027) | 0.805 | 0.816 | 0.615 |
| *R+S*(*EX*) [22](2022,TCSVT) | 0.836(0.031) | 0.831(0.025) | 0.641(0.032) | 0.832(0.023) | 0.833(0.019) | 0.634(0.017) | 0.833 | 0.832 | 0.636 |
| **DisCoVQA (Ours)** | 0.820(0.030) | 0.826(0.024) | 0.633(0.034) | 0.847(0.014) | 0.847(0.028) | 0.660(0.018) | 0.838 | 0.840 | 0.651 |

TABLE III: SROCC comparison of different methods with 8:2 setting and 6:2:2 setting on LIVE-VQC and KoNViD-1k.

| Dataset / | LIVE-VQC | | KoNViD-1k | |
|---|---|---|---|---|
| Method / Split Setting | 6:2:2 | 8:2 | 6:2:2 | 8:2 |
| CNN-TLVQM(*EX*) [21] | 0.830 | 0.814 | 0.830 | 0.816 |
| PVQ(*EX*) [42] | 0.827 | 0.803 | 0.791 | 0.795 |
| LSCT-PHIQ(*EX*) [43] | NA | 0.796 | 0.860 | 0.833 |
| **DisCoVQA (Ours)** | 0.838 | 0.820 | 0.863 | 0.847 |

TABLE IV: PLCC comparison of different methods with 8:2 setting and 6:2:2 setting on LIVE-VQC and KoNViD-1k.

| Dataset / | LIVE-VQC | | KoNViD-1k | |
|---|---|---|---|---|
| Method / Split Setting | 6:2:2 | 8:2 | 6:2:2 | 8:2 |
| CNN-TLVQM(*EX*) [21] | 0.840 | 0.821 | 0.830 | 0.818 |
| PVQ(*EX*) [42] | 0.837 | 0.811 | 0.786 | 0.774 |
| LSCT-PHIQ(*EX*) [43] | NA | 0.782 | 0.850 | 0.834 |
| **DisCoVQA (Ours)** | 0.844 | 0.826 | 0.860 | 0.847 |

discussed above. We follow the standard 6:2:2 train-validate-test dataset split settings (60% for training, 20% for validation, and we report our performance on *the rest 20% test set* while the validation performance reaches peak), and report the average results on ten random splits for each dataset, together with the standard deviations. This evaluation setting is to avoid overfitted results and improve the confidence of our experimental setting.

### B. Comparison on Individual Datasets

As Tab. II shows, our proposed DisCoVQA has shown superior performance than existing methods published prior

to us on four individual benchmark datasets, even those with extra-dataset pretraining (labeled as *EX*). For example, the proposed model is up to **9.5%** better than VSFA [23] with similar parameters and computational complexities. The proposed model is also up to **13.8%** better than MLSP-FF which has much more parameters but no temporal modeling other than average pooling, showing the vitality of proper temporal modeling in VQA.

The proposed model is **state-of-the-art on three datasets** and the runner-up model on LIVE-VQC, slightly inferior to *R+S* [22]. However, *R+S* includes an additional spatial branch that is fine-tuned on other IQA datasets. The proposed model is

TABLE V: Cross-dataset results: between LIVE-VQC [33], KoNViD-1k [17] and YouTube-UGC [41]. Results of DisCoVQA related to YouTube-UGC [41] is conducted with the 236-fewer-video version. We still reach better cross-dataset performance when training on YouTube-UGC with less videos (our results related to YouTube-UGC are labeled with $\star$ and presented only for reference due to video missing.)

| Train on | KoNViD-1k | | | | LIVE-VQC | | | | Youtube-UGC | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Test on | LIVE-VQC | | Youtube-UGC | | KoNViD-1k | | Youtube-UGC | | LIVE-VQC | | KoNViD-1k | |
| | SROCC | PLCC | SROCC | PLCC | SROCC | PLCC | SROCC | PLCC | SROCC | PLCC | SROCC | PLCC |
| CNN-TLVQM [21] | 0.713 | 0.752 | NA | NA | 0.642 | 0.631 | NA | NA | NA | NA | NA | NA |
| GST-VQA [5] | 0.700 | 0.733 | NA | NA | 0.709 | 0.707 | NA | NA | NA | NA | NA | NA |
| VIDEVAL [38] | 0.627 | 0.654 | 0.370 | 0.390 | 0.625 | 0.621 | 0.302 | 0.318 | 0.542 | 0.553 | 0.610 | 0.620 |
| MDTVSFA [24] | 0.716 | 0.759 | 0.408 | 0.443 | 0.706 | 0.711 | 0.355 | 0.388 | 0.582 | 0.603 | 0.649 | 0.646 |
| **DisCoVQA (Ours)** | 0.782 | 0.797 | 0.415$\star$ | 0.449$\star$ | 0.792 | 0.785 | 0.409$\star$ | 0.432$\star$ | 0.661$\star$ | 0.685$\star$ | 0.686$\star$ | 0.697$\star$ |

**14.5%**[3] better than its pure-temporal branch and **5.7%** better than it while their spatial branch is not extra-pretrained. This comparison suggests that at least the proposed transformer-based approach is competitive for temporal modeling. [4]

We also notice different datasets, though all aimed at collecting natural distortions, have different characteristics, especially considered in the temporal domain. For example, the KoNViD-1k contains more **diverse contents across time** and **compression-based artifacts**, while LIVE-VQC contains more **in-capture temporal distortions**. Therefore, the *weighted average* performance on these two datasets (Tab. II, in the rightmost) might be a more reliable benchmark to evaluate the full ability of methods. The proposed DisCoVQA reaches state-of-the-art on group weighted averages for both simulated and real-world natural datasets, showing that the proposed method is generally robust rather specially effective on some specific distortion types.

We also notice that some recent deep methods [21], [42], [43] directly report the performance on the validation set in a different 8:2 setting (80% for training, 20% for validation, no extra testing dataset). To evaluate the difference of our setting and this setting, we also evaluated the proposed DisCoVQA on this setting and try our best to reproduce them in our setting. The results in Tab. IV and Tab. III suggested that two different settings have statistically prominent performance gap, indicating that we might need to align these settings before making fair comparisons between different methods.

### C. Cross-dataset Results

To measure the generalization ability of the proposed DisCoVQA, we compare cross-dataset results with several state-of-the-art methods. We carefully choose the methods with relatively good generalization ability for comparison: MDTVSFA [24] is specially designed for multi-dataset alignment; GST-VQA [5] is also designed for better generalization; CNN-TLVQM [21], CoINVQ [40] and VIDEVAL [38] ensemble different types of features for robustness. Without including any special design, the proposed DisCoVQA reaches better generalization than them. The cross-dataset results among

KoNViD-1k [17], LIVE-VQC [33] and YouTube-UGC [41][5] are reported in Tab. V.

Compared with existing methods, we observe more obvious improvements in cross-dataset experiments. Take our comparison with CNN-TLVQM [21] as an example. CNN-TLVQM only relies on handcraft features for temporal modeling, while we design both temporal distortion and temporal attention modeling with transformer-based backbones for it. As Tab. II and Tab. V shows, we outperform CNN-TLVQM by **10%** in cross-dataset results (between LIVE-VQC & KoNViD-1k) where the proposed model is only 2% better than it during intra-dataset settings. It demonstrates that though applying handcraft or other traditional solutions to tackle the temporal relationships can reach good performance, the proposed transformer-based approach can learn more these relationships more robustly.

We also notice that current methods still cannot generalize well between KoNViD-1k/LIVE-VQC and YouTube-UGC, which might be due to the large proportion of non-photographic videos (categories *games*, *animation*, *lyric videos*, *news report* in it) in YouTube-UGC, while there are very few of them in LIVE-VQC and KoNViD-1k. It will be a nice future objective to improve the generalization ability of VQA approaches between these generated videos and other natural videos.

### D. Results on Large-scale Datasets

We evaluate the proposed model on two large-scale datasets, LSVQ [42] and KoNViD-150k [13]. As Tab. VI shows, we outperform the PVQ [42] with **7.7%** improvement on the same setting, and 5.1% improvement even when PVQ uses the additional '*patch*' annotations. The advanced performance of DisCoVQA has shown the effectiveness of the proposed temporal modeling methodology. We also outperform the only available model trained on KoNViD-150k, the MLSP-FF [13] which proposed the dataset and VSFA (reproduced by us).

We further notice that the cross-dataset performance on KoNViD-1k [17] and LIVE-VQC [33] of DisCoVQA trained with LSVQ is obviously better than that of KoNViD-150k. We suspect that this might be due to the different temporal

---

[3]Result based on their paper. Same for the next.

[4]Some methods did not provide their codes or report their full performance, so we directly reported results in their paper and left their missing results empty. We try our best to reproduce and fill in every comparison for every competitive method.

[5]The $\star$-*labeled* results related to YouTube-UGC are shown for reference only due to difference of dataset size (our used are 236 fewer than the original version due to missing of the downloadable links).

TABLE VI: Large Dataset Training I: pretrained with large-scale LSVQ [42] and evaluate on different sets without fine-tuning.

| Train on | $LSVQ_{train}$ | | | | | |
|---|---|---|---|---|---|---|
| Test on | LIVE-VQC | | KoNViD-1k | | $LSVQ_{test}$ | |
| | SROCC | PLCC | SROCC | PLCC | SROCC | PLCC |
| VSFA [23] | 0.734 | 0.772 | 0.784 | 0.794 | 0.801 | 0.796 |
| TLVQM [20] | 0.670 | 0.691 | 0.732 | 0.724 | 0.772 | 0.774 |
| PVQ [42] | 0.747 | 0.776 | 0.781 | 0.781 | 0.814 | 0.816 |
| *PVQ+* | *0.770* | *0.807* | *0.791* | *0.795* | *0.827* | *0.828* |
| **DisCoVQA (Ours)** | 0.823 | 0.837 | 0.846 | 0.849 | 0.859 | 0.850 |

TABLE VII: Large Dataset Training II: pretrained with large-scale KoNViD-150k [13] and evaluate on different sets without fine-tuning.

| Train on | KoNViD-150k-A | | | | | |
|---|---|---|---|---|---|---|
| Test on | LIVE-VQC | | KoNViD-1k | | KoNViD-150k-B | |
| | SROCC | PLCC | SROCC | PLCC | SROCC | PLCC |
| VSFA [23] | 0.708 | 0.733 | 0.801 | 0.815 | 0.813 | 0.808 |
| MLSP-FF [13] | 0.738 | 0.754 | 0.828 | 0.821 | 0.827 | 0.852 |
| **DisCoVQA (Ours)** | 0.751 | 0.766 | 0.843 | 0.841 | 0.845 | 0.858 |



(a) Examples in KoNViD-150k.
Less complicated temporal relationships.

(b) Examples in LSVQ.
More complicated temporal relationships.

Fig. 4: Comparison of (a) KoNViD-150k [13] and (b) LSVQ [42]. The LSVQ dataset has shown more complicated temporal relationships (more difference between frame contents; more temporal distortions) than the KoNViD-150k.

TABLE VIII: Running time comparison for DisCoVQA on a 540P, 8-second video in KoNViD-1k [17] dataset. The standard deviants for the corresponding running times are in brackets.

| **Backbone Extractor** / | ResNet-50 | I3D-ResNet-50 | Swin-T (as proposed) |
|---|---|---|---|
| Running Time | 2.71s(0.14s) | 2.45s(0.17s) | **2.33s(0.09s)** |
| **Temporal Regression** / | 6-layer GRU | vanilla Transformer | the TCT (proposed) |
| Running Time | 0.053s(0.011s) | 0.120s(0.018s) | **0.021s(0.004s)** |

relationships in two datasets. As illustrated in Fig. 4, the temporal relationships (both temporal distortions and different contents of different frames) are less complicated in KoNViD-150k than in LSVQ, which might be caused by the different collection sources of these two datasets and the different average durations (KoNViD-150k: 5s, LSVQ: $\tilde{}7.5$s). The LSVQ dataset with more complicated temporal relationships have reached better generalization ability in natural VQA datasets, suggesting that temporal issues are common in natural videos.

### E. Running Time Comparison

We discuss the computational cost introduced by transformers in two parts: the running time for the transformer-based backbone (Swin-T) in the STDE and for the transformer-based TCT, and compare with several existing alternative approaches (CNN&RNN), as shown in Tab. VIII. As the results shows, switching the traditional CNN/3D-CNN backbones into the transformer-based Swin-T backbone does not lead to additional running time. The proposed TCT with the temporal sampling on features (TSF) is also **6x** faster than vanilla transformer and **2.5x** faster than GRU, proving that it alleviates the problem of high computation loads of transformers on long sequences. These results prove that the proposed DisCoVQA maintains the efficiency of existing deep VQA models while reaching remarkably better performance with transformer-based architectures.

### V. EXPERIMENTAL ANALYSIS ON MODEL DESIGN

In this part, we would like to answer three important questions to evaluate the effectiveness of the proposed model:

1) Does every design of the proposed model lead to reasonable performance improvement? (In Sec. V-A and Sec. V-B)

2) In which particular cases can the model show better ability? (In Sec. V-C)

3) Does the model give reliable results as the TSF includes some randomness during inference? (In Sec. V-D.)

Without special notes, the training datasets for all these studies are LSVQ dataset due to its large scale (28K training videos) and diversity on temporal relationships (as discussed in Sec. IV-D).

### A. Ablation Studies on the STDE

To discuss the effectiveness of the proposed STDE, we run ablation studies on LSVQ dataset [42], and provide result comparisons on both intra-dataset and cross-dataset (test on LIVE-VQC and KoNViD-1k) experiments. We run ablation studies for both the backbone network and micro-designs (multi-level extraction & temporal differences) in the STDE.

*1) Effects of Backbone Structures:* In Tab. IX, we replace the Swin-T backbone into two CNN backbones with similar parameters and running time (as compared in Sec. IV-E): the ResNet-50 that does not extract temporal quality information, and the I3D-ResNet-50 that extracts temporal quality information with convolution kernels. The I3D-ResNet-50 backbone is better than ResNet-50 backbone as it has temporal sensitivity, but ours with Swin-T still performs notably better than it, especially on LIVE-VQC test set with most severe temporal distortions. It should also be noted that even with the same backbone (ResNet50), our method still has better performance than VSFA, proving that our designs other than the backbone are also effective.

*2) Effects of Micro-designs: Multi-level & Temporal Differences:* We discuss about the two important designs in the STDE: the multi-level feature extraction and the temporal differences. As shown in Tab. X, the multi-level extraction

TABLE IX: Ablation study on backbone structures in the STDE. All these backbones have similar parameters and running time for feature extraction (Tab. VIII).

| Testing on / | LSVQ$_{test}$ | | KoNViD-1k | | LIVE-VQC | |
|---|---|---|---|---|---|---|
| Backbone Network | SROCC | PLCC | SROCC | PLCC | SROCC | PLCC |
| VSFA [23] (with ResNet-50) | 0.801 | 0.796 | 0.784 | 0.794 | 0.734 | 0.772 |
| ResNet-50 [15] | 0.823 | 0.822 | 0.810 | 0.815 | 0.758 | 0.776 |
| I3D-ResNet-50 [4] | 0.840 | 0.832 | 0.825 | 0.817 | 0.774 | 0.793 |
| **Swin-T** (as proposed) | **0.859** | **0.850** | **0.846** | **0.849** | **0.823** | **0.837** |

TABLE X: Ablation study on multi-level extraction and temporal differences. Both micro-designs further improves the performance of transformer-based STDE.

| Testing on / | LSVQ$_{test}$ | | KoNViD-1k | | LIVE-VQC | |
|---|---|---|---|---|---|---|
| Micro-Design | SROCC | PLCC | SROCC | PLCC | SROCC | PLCC |
| *w/o* Both (remove micro-designs) | 0.828 | 0.824 | 0.819 | 0.824 | 0.786 | 0.798 |
| *w/o* Multi-level Extraction | 0.835 | 0.829 | 0.825 | 0.828 | 0.794 | 0.809 |
| *w/o* Temporal Differences | 0.844 | 0.835 | 0.836 | 0.839 | 0.806 | 0.822 |
| **Full STDE** (proposed) | **0.859** | **0.850** | **0.846** | **0.849** | **0.823** | **0.837** |

significantly improves the accuracy on all testing scenarios, supporting our claim that it enhances the distortion sensitivity of the STDE. The temporal differences also show non-trivial improvements on all sets, and is especially helpful (+2%) on LIVE-VQC test set with all hand-held video shots, where the temporal distortions are most severe. The combination of these two designs lead to 3% improvement than the vanilla Swin-T extractor (*w/o* Both in the Table). Both micro-designs help the proposed transformer-based STDE to be more suitable for extracting temporal distortions.

### B. Ablation Studies on the TCT

In this part, we compare the proposed TCT with two groups of variants on LSVQ dataset. The first group of variants are non-transformer structures: the temporal MLP, the temporal CNN, and the LSTM (a type of RNN); the second group of variants are the structural variants of transformers. We also compare the model variant that removes the TCT at all. Moreover, we discuss the effects for the temporal sampling on features (TSF) for training the TCT on different scale of datasets to show how they help to improve the hard cases of introducing transformers in VQA.

*1) Comparison with Non-Transformer Structures:* In this part, we compare the transformer-based TCT with several non-transformer structures. We set these structures with the same layers (6 layers) as the TCT to make fair comparisons and the results are shown in Tab. XI, Group 1. The proposed transformer-based TCT has much better performance than all non-transformer variants and the variant that removes the TCT, demonstrating that transformer architectures are better for temporal quality attention modeling in VQA.

*2) Comparison with Transformer-based Variants:* In this part, we compare the proposed encoder-decoder-like TCT with several variants, including the pure transformer encoder with 4/6 layers and the variant that changes the average content token ($\mathbf{T}_{avg}$) into a zero token as the target input of the

TABLE XI: Ablation study of the TCT architecture: compared with non-transformer structures and structural variants of transformers. The corresponding running time comparison can be found in Tab. VIII.

| Testing on / | LSVQ$_{test}$ | | KoNViD-1k | | LIVE-VQC | |
|---|---|---|---|---|---|---|
| Variant | SROCC | PLCC | SROCC | PLCC | SROCC | PLCC |
| *remove the TCT* | 0.842 | 0.836 | 0.826 | 0.823 | 0.804 | 0.813 |
| Group 1: Non-Transformer Structures | | | | | | |
| Temporal MLP (6-layer) | 0.836 | 0.828 | 0.824 | 0.818 | 0.792 | 0.801 |
| Temporal CNN (6-layer) | 0.839 | 0.832 | 0.828 | 0.822 | 0.799 | 0.808 |
| LSTM (RNN, 6-layer) | 0.841 | 0.838 | 0.831 | 0.830 | 0.806 | 0.816 |
| Group 2: Structural Variants of Transformers | | | | | | |
| Pure Encoder (4-layer) | 0.847 | 0.840 | 0.837 | 0.838 | 0.812 | 0.826 |
| Pure Encoder (6-layer) | 0.848 | 0.842 | 0.841 | 0.842 | 0.812 | 0.828 |
| *change* $\mathbf{T}_{avg}$ *as zero token* | 0.852 | 0.845 | 0.843 | 0.844 | 0.813 | 0.832 |
| **Full TCT** (proposed) | **0.859** | **0.850** | **0.846** | **0.849** | **0.823** | **0.837** |

TABLE XII: Ablation study of the proposed temporal sampling on features (TSF) on small and long-duration VQA datasets.

| Dataset / | CVD2014 | | LIVE-Qualcomm | | YouTube-UGC | |
|---|---|---|---|---|---|---|
| **Size/Average Duration** | 234/10s | | 208/15s | | 1144/20s | |
| Strategy | SROCC | PLCC | SROCC | PLCC | SROCC | PLCC |
| full-length features (*w/o* TSF) | 0.878 | 0.879 | 0.803 | 0.805 | 0.789 | 0.790 |
| pooling on features | 0.883 | 0.880 | 0.799 | 0.801 | 0.794 | 0.797 |
| sampling on features (proposed) | **0.897** | **0.893** | **0.823** | **0.825** | **0.809** | **0.808** |
| *improvement* to *w/o* TSF | **+2.2%** | **+1.6%** | **+2.6%** | **+2.9%** | **+2.5%** | **+2.3%** |

transformer decoder. Changing the proposed structure into all these variants result in **2%** performance drop, proving the effectiveness of the proposed encoder-decoder-like structure of the TCT which takes the average content as target input. This result also suggests that the overall content is vital in deciding the temporal quality attention across frames in VQA.

*3) Effects of temporal sampling on features (TSF):* We show the effectiveness of implementing the TSF for the TCT regression in several different datasets. First, as shown in Tab. XII, for datasets that are either small (LIVE-Qualcomm and CVD2014) or with long duration (YouTube-UGC), the TSF significantly helps the learning process of the TCT. It is also noteworthy that on other datasets (as compared in Tab. XIII), the TSF does not lead to noticeable better performance. As TSF also significantly improves the training speed, we still take it as a part of the TCT when training DisCoVQA in these datasets. Also, the temporal pooling on features consistently perform worse than the TSF on six datasets, proving that keeping original features is important

TABLE XIII: Ablation study of the TSF on other datasets.

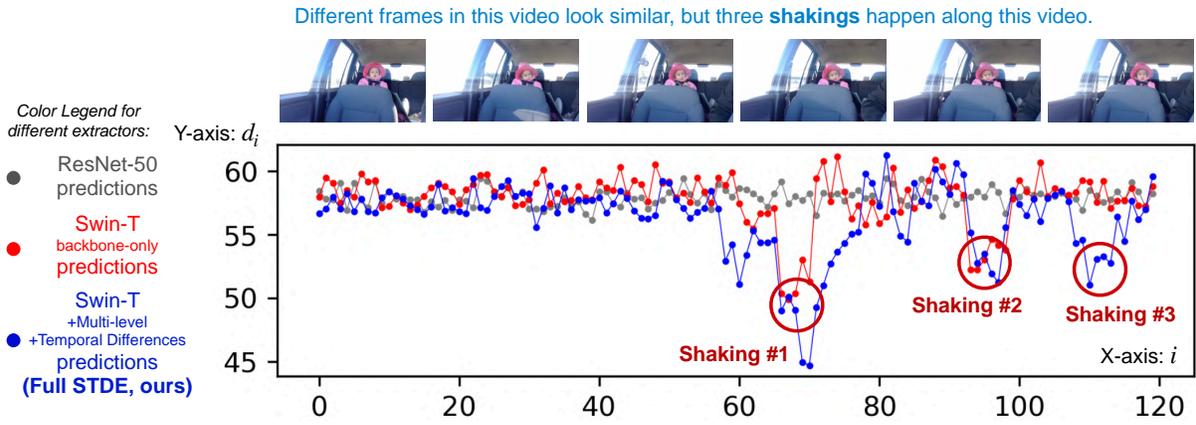| Dataset / | LIVE-VQC | | KoNViD-1k | | LSVQ | |
|---|---|---|---|---|---|---|
| **Size/Average Duration** | 585/8s | | 1200/10s | | 39075/7s | |
| Strategy | SROCC | PLCC | SROCC | PLCC | SROCC | PLCC |
| full-length features (w/o TSF) | 0.816 | 0.823 | 0.846 | 0.847 | **0.860** | **0.850** |
| pooling on features | 0.809 | 0.818 | 0.839 | 0.840 | 0.857 | 0.846 |
| sampling of features (proposed) | **0.820** | **0.826** | **0.847** | **0.847** | 0.859 | **0.850** |
| *improvement* to *w/o* TSF | **+0.4%** | **+0.3%** | **+0.1%** | 0.0% | -0.1% | 0.0% |

Fig. 5: Visualizations of temporal-distortion-aware qualities $d_i$ learnt in STDE, compared with different variants of STDE. Discussions are in Sec. V-C.
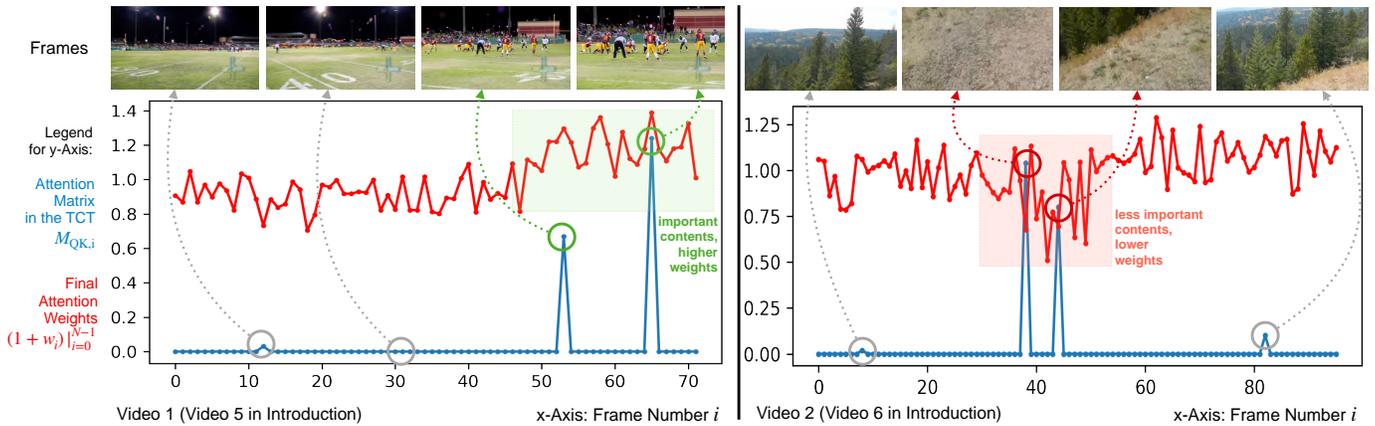


Fig. 6: Visualizations of temporal quality attention, including the attention matrix in TCT and the final temporal quality attention weights, demonstrating that the propose TCT can learn reasonable temporal quality attention. Discussions are in Sec. V-C.

to the learning of TCT. The TSF provides a feasible way of implementing transformers into VQA both effectively and efficiently.

### C. Qualitative Results

*1) Visualization of Temporal Distortion Sensitivity:* For the STDE, we aim to spot the temporal distortions between frames, such as shaking and flickers in videos. In Fig. 5, we visualize the distortion-only qualities of frames in a video that contain shaking learnt with features extracted by different approaches. As illustrated in the figure, compared to vanilla feature extraction from CNN-based backnones, the adoption of Swin-T backbone and two micro-designs in STDE both enhances the temporal distortion extraction by detecting more shakings, and helps to better spot temporal distortions in STDE.

*2) Visualization of Temporal Quality Attention:* For the TCT, we aim to enhance the weights of important frames (frames more close to theme) and suppress the weights of unimportant frames (irrelevant frames). In Fig. 6, we visualize the attention matrix $M_{QK}$ and final attention weights $\{w_i|_{i=0}^{N-1}\}$ learnt in the TCT. For video 1, the peak attention

is at the zoom-in for players, which are the specially important frames for this replay video of the football match, and they result in higher weights for these frames. For video 2, on the contrary, the peak attention comes at the irrelevant frames that photographs on the ground, which are specially irrelevant frames and have lower final weights. These results demonstrates the effectiveness of the proposed TCT.

*3) Visualization for Correlations with Ground Truth:* To further visualize the result of the proposed DisCoVQA, we show the correlation of DisCoVQA predicted scores and the ground truth labels in Fig. 7. The x-axis represents ground truth labels (MOS), and the y-axis shows the predicted scores $\hat{q}$ with respect to MOS. The bright blue line is the reference line when the prediction is the same as MOS. In visualizing these correlations, we add two gray lines as **unit deviation** for each dataset (the quality prediction is roughly correct if fallen in this line). We find out that the proposed DisCoVQA consistently predicts video scores with a very high correlation with ground truth scores, and only a few videos fall out of the range enclosed by the gray lines. These correlations, together as results shown in Tab. II, demonstrate that the proposed DisCoVQA is an accurate quality evaluator.
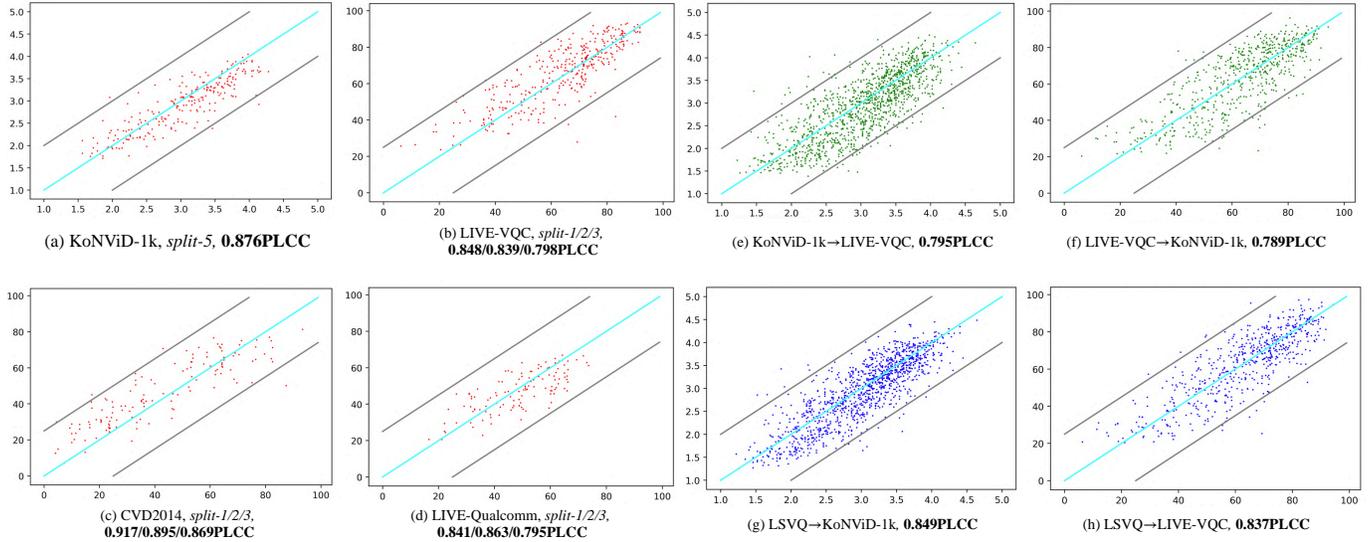
Fig. 7: The correlation of predicted scores (**y-axis**) and the ground truth labels (**x-axis**). (a)(b)(c)(d) show settings within individual datasets, and (e)(f)(g)(h) show cross-dataset settings. Corresponding benchmarks are in Tab. II, Tab. V and Tab. VI.

TABLE XIV: Reliability analysis for TSF I: the metric results with respect to sample count $M$, with train set LSVQ$_{train}$ and test set LIVE-VQC.

| Sample Count $S_m$ | 1 | | 2 | | 4 | | 8 | | $\infty$(40) | |
|---|---|---|---|---|---|---|---|---|---|---|
| | SROCC | PLCC | SROCC | PLCC | SROCC | PLCC | SROCC | PLCC | SROCC | PLCC |
| **DisCoVQA** (proposed) | 0.816 | 0.824 | 0.819 | 0.831 | 0.821 | 0.836 | **0.823** | **0.837** | **0.823** | **0.837** |

TABLE XV: Reliability analysis for TSF II: the mean standard deviants $\sigma_M$ of different predictions on the same video, with train set LSVQ$_{train}$ and test set LIVE-VQC.

| Sample Count $S_m$ | 1 | 2 | 4 | 6 | 8 | 16 | 40 |
|---|---|---|---|---|---|---|---|
| **DisCoVQA** (proposed) | 0.0114 | 0.0043 | 0.0025 | 0.0016 | **0.0010** | **0.0006** | **0.0002** |

*4) Visualization for Success and Failure Cases:* We also visualize several successful or failed prediction of the proposed model in LSVQ dataset. As the plot in Fig. 8(d) shows, the proposed DisCoVQA can give reasonable quality predictions on most videos, including the video Fig. 8(c) with complex changing contents across the video. We also visualize two specific failure cases of the model: (a) **non-natural contents**, which have been discussed above and are also hard situations for several existing methods such as PVQ [42]; (b) **ambiguity of human annotations** (this video contains strong flicker but still gets relatively high MOS), which suggests that the proposed model is especially sensitive on temporal distortions, though for this case the human annotators prefer to give it higher MOS scores. Results of these cases are in line with our analysis for the proposed method.

### D. Reliability Analysis for TSF during inference

We have proposed the TSF to significantly reduce the computational complexity of transformer-based modules in VQA. During inference, we random sample features for $S_m = 8$ times and get the average quality prediction for these samples.
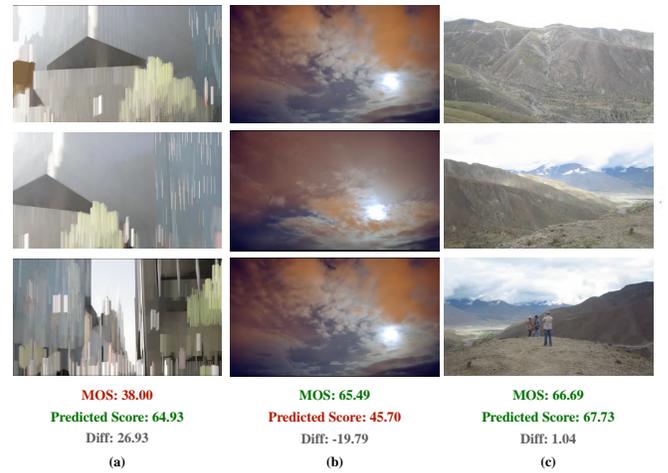


Fig. 8: The success and failure cases on dataset LSVQ [42] by the proposed DisCoVQA. (a) **non-natural contents** and (b) **ambiguity on human annotations** (on a flicker video) are two representative failure cases, where (c) is one of the successful predictions.

However, as the TSF contains random sampling for frames in each segment, we need to do the following two experiments to confirm its reliability. First, we need to confirm that the TSF will not reduce the overall prediction accuracy of the proposed DisCoVQA. We choose to evaluate on the setting with training set LSVQ$_{train}$ and testing set LIVE-VQC, and as the results in Tab. XIV demonstrates, we can obtain nearly the same result as $S_m = 8$ compared with infinite samples ($S_m = 40$ in practice), proving that the TSF does not harm the overall prediction accuracy. Moreover, we show the mean standard deviants $\sigma_m$ of predictions on different samples (normalized with the overall score range) of the same video to $S_m$ in

Tab. XV, proving that with $S_m \geq 8$, the $\sigma_m$ can be negligible (less than 0.001), thus DisCoVQA with TSF can still infer with high stability. Both experiments prove that TSF is reliable during inference.

## VI. CONCLUSION AND FUTURE WORKS

We have proposed **DisCoVQA**, a novel and effective method that aims at better modeling both temporal distortions and content-related temporal quality attention via transformer-based architectures. To better capture temporal distortions, we extract multi-level features from a Swin-T backbone network for a better semantic understanding of video actions and compute temporal differences to further spot temporal variations. To model the temporal quality attention toward different importance of frames, we utilize a transformer encoder-decoder structure to consider the correlation of frame contents to the overall video theme. We also introduce the temporal sampling on features to boost the training effectiveness and efficiency of this transformer-based temporal regression module. In conclusion, we propose a transformer-based method that better models the temporal relationships in VQA, and the proposed DisCoVQA has reached state-of-the-art performance on several natural VQA datasets and achieved excellent generalization ability among them.

In the future, we aim at solving several problems not well addressed by current frameworks (as analyzed in failure cases in Fig. 8), including the better coverage of non-natural contents, and dealing with ambiguous quality scores. We also notice that several recent methods benefit from extra pre-training, yet they all need labeled datasets. For the next step, we hope to propose a method to include label-free pre-training for VQA that can lead to further improvements on performance.

## REFERENCES

[1] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lucic, and Cordelia Schmid. Vivit: A video vision transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6836–6846, October 2021.

[2] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding?, 2021.

[3] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. 2020.

[4] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.

[5] Baoliang Chen, Lingyu Zhu, Guo Li, Fangbo Lu, Hongfei Fan, and Shiqi Wang. Learning generalized spatial-temporal deep feature representation for no-reference video quality assessment. *IEEE Transactions on Circuits and Systems for Video Technology*, 2021.

[6] Kyunghyun Cho, Bart van Merrienboer, Çaglar Gülçehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder-decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1724–1734. ACL, 2014.

[7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics.

[8] Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra Malik, and Christoph Feichtenhofer. Multiscale vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6824–6835, October 2021.

[9] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6201–6210, 2019.

[10] Karl Friston. Friston, k.j.: The free-energy principle: a unified brain theory? nat. rev. neurosci. 11, 127-138. *Nature reviews. Neuroscience*, 11:127–38, 02 2010.

[11] Deepti Ghadiyaram and Alan C. Bovik. Perceptual quality prediction on authentically distorted images using a bag of features approach. *Journal of Vision*, 17(1):32–32, 01 2017.

[12] Deepti Ghadiyaram, Janice Pan, Alan C. Bovik, Anush Krishna Moorthy, Prasanjit Panda, and Kai-Chieh Yang. In-capture mobile video distortions: A study of subjective behavior and objective algorithms. *IEEE Transactions on Circuits and Systems for Video Technology*, 28(9):2061–2077, 2018.

[13] Franz Götz-Hahn, Vlad Hosu, Hanhe Lin, and Dietmar Saupe. Konvid-150k: A dataset for no-reference video quality assessment of videos in-the-wild. In *IEEE Access 9*, pages 72139–72160. IEEE, 2021.

[14] Kensho Hara, Hirokatsu Kataoka, and Yutaka Satoh. Learning spatio-temporal features with 3d residual networks for action recognition. In *2017 IEEE International Conference on Computer Vision Workshops (ICCVW)*, pages 3154–3160, 2017.

[15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.

[16] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.

[17] Vlad Hosu, Franz Hahn, Mohsen Jenadeleh, Hanhe Lin, Hui Men, Tamás Szirányi, Shujun Li, and Dietmar Saupe. The konstanz natural video database (konvid-1k). In *Ninth International Conference on Quality of Multimedia Experience (QoMEX)*, pages 1–6, 2017.

[18] Junjie Ke, Qifei Wang, Yilin Wang, Peyman Milanfar, and Feng Yang. Musiq: Multi-scale image quality transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5148–5157, October 2021.

[19] Woojae Kim, Jongyoo Kim, Sewoong Ahn, Jinwoo Kim, and Sanghoon Lee. Deep video quality assessor: From spatio-temporal visual sensitivity to a convolutional neural aggregation network. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.

[20] Jari Korhonen. Two-level approach for no-reference consumer video quality assessment. *IEEE Transactions on Image Processing*, 28(12):5923–5938, 2019.

[21] Jari Korhonen, Yicheng Su, and Junyong You. Blind natural video quality prediction via statistical temporal features and deep spatial features. In *Proceedings of the 28th ACM International Conference on Multimedia*, MM '20, page 3311–3319, New York, NY, USA, 2020. Association for Computing Machinery.

[22] Bowen Li, Weixia Zhang, Meng Tian, Guangtao Zhai, and Xianpei Wang. Blindly assess quality of in-the-wild videos via quality-aware pre-training and motion perception. *IEEE Transactions on Circuits and Systems for Video Technology*, 2022.

[23] Dingquan Li, Tingting Jiang, and Ming Jiang. Quality assessment of in-the-wild videos. In *Proceedings of the 27th ACM International Conference on Multimedia*, MM '19, page 2351–2359, New York, NY, USA, 2019. Association for Computing Machinery.

[24] Dingquan Li, Tingting Jiang, and Ming Jiang. Unified quality assessment of in-the-wild videos with mixed datasets training. *International Journal of Computer Vision*, 129(4):1238–1257, 2021.

[25] Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video swin transformer. *arXiv preprint arXiv:2106.13230*, 2021.

[26] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019.

[27] Anish Mittal, Anush Krishna Moorthy, and Alan Conrad Bovik. No-reference image quality assessment in the spatial domain. *IEEE Transactions on Image Processing*, 21(12):4695–4708, 2012.

[28] Anish Mittal, Michele A. Saad, and Alan C. Bovik. A completely blind video integrity oracle. *IEEE Transactions on Image Processing*,

25(1):289–300, 2016.

[29] Mikko Nuutinen, Toni Virtanen, Mikko Vaahteranoksa, Tero Vuori, Pirkko Oittinen, and Jukka Häkkinen. Cvd2014—a database for evaluating no-reference video quality assessment algorithms. *IEEE Transactions on Image Processing*, 25(7):3073–3086, 2016.

[30] Mariusz Oszust. No-reference image quality assessment with local features and high-order derivatives. *Journal of Visual Communication and Image Representation*, 56:15–26, 2018.

[31] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.

[32] Michele A. Saad, Alan C. Bovik, and Christophe Charrier. Blind image quality assessment: A natural scene statistics approach in the dct domain. *IEEE Transactions on Image Processing*, 21(8):3339–3352, 2012.

[33] Zeina Sinno and Alan Conrad Bovik. Large-scale study of perceptual video quality. *IEEE Transactions on Image Processing*, 28(2):612–627, 2019.

[34] Shaolin Su, Qingsen Yan, Yu Zhu, Cheng Zhang, Xin Ge, Jinqiu Sun, and Yanning Zhang. Blindly assess image quality in the wild guided by a self-adaptive hyper network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

[35] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alexander A. Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*, AAAI'17, page 4278–4284. AAAI Press, 2017.

[36] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning spatiotemporal features with 3d convolutional networks. In *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 4489–4497, 2015.

[37] Zhengzhong Tu, Chia-Ju Chen, Yilin Wang, Neil Birkbeck, Balu Adsumilli, and Alan C. Bovik. Efficient user-generated video quality prediction. In *2021 Picture Coding Symposium (PCS)*, pages 1–5, 2021.

[38] Zhengzhong Tu, Yilin Wang, Neil Birkbeck, Balu Adsumilli, and Alan C. Bovik. Ugc-vqa: Benchmarking blind video quality assessment for user generated content. *IEEE Transactions on Image Processing*, 30:4449–4464, 2021.

[39] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS'17, page 6000–6010, Red Hook, NY, USA, 2017. Curran Associates Inc.

[40] Yilin Wang, Junjie Ke, Hossein Talebi, Joong Gon Yim, Neil Birkbeck, Balu Adsumilli, Peyman Milanfar, and Feng Yang. Rich features for perceptual quality assessment of ugc videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13435–13444, June 2021.

[41] Joong Gon Yim, Yilin Wang, Neil Birkbeck, and Balu Adsumilli. Subjective quality assessment for youtube ugc dataset. In *2020 IEEE International Conference on Image Processing (ICIP)*, pages 131–135, 2020.

[42] Zhenqiang Ying, Maniratnam Mandal, Deepti Ghadiyaram, and Alan Bovik. Patch-vq: 'patching up' the video quality problem. In *2021 Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14019–14029, June 2021.

[43] Junyong You. Long short-term convolutional transformer for no-reference video quality assessment. In *Proceedings of the 29th ACM International Conference on Multimedia*, MM '21, page 2112–2120, New York, NY, USA, 2021. Association for Computing Machinery.

[44] Junyong You and Jari Korhonen. Deep neural networks for no-reference video quality assessment. In *Proceedings of the IEEE International Conference on Image Processing (ICIP)*, pages 2349–2353, 2019.