# CAMP-VQA: Caption-Embedded Multimodal Perception for No-Reference Quality Assessment of Compressed Video
# Supplementary Material for WACV 2026

Xinyi Wang     Angeliki Katsenou     Junxiao Shen     David Bull

School of Computer Science, University of Bristol

Bristol, United Kingdom

`{xinyi.wang, angeliki.katsenou, junxiao.shen, dave.bull}@bristol.ac.uk`

## Abstract

*This document is the supplementary material for our WACV submission. We provide further explanation and details on frame difference fragmentation, quality-aware caption generation, spatio-temporal feature extraction, and loss functions. Additional figures are included to illustrate the quality-related prompt settings and to demonstrate the effectiveness of the proposed CAMP-VQA in predicting scores across different quality dimensions.*
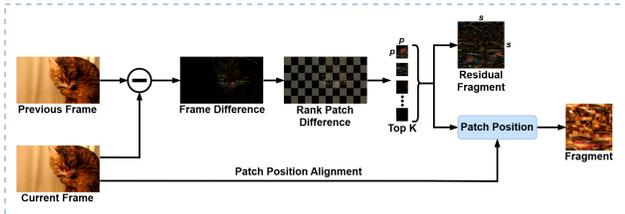
## A. Frame Difference Fragmentation



Figure 1. Frame difference fragmentation (FDF) module.

Regarding patch selection, we compute the inter-frame absolute differences for all patches and rank them accordingly. The fragmentation details are shown in (Fig. 1). The number of selected patches $K$ is calculated based on the target model input size $s \times s$ and each patch size $p \times p$:

$$K = \frac{s^2}{p^2}. \tag{1}$$

We then select the top $K$ patches from the ranked list based on the sum of absolute inter-frame differences. To ensure spatial consistency between the extracted frame and residual fragments, each selected patch in the residual is paired with a corresponding fragment extracted from the same pixel position in $F_t$. For a target size of $s = 224$, we use a

patch size of $p = 16$, resulting in $K = 196$. This method generalizes across resolutions, as fragmentation operates directly on the original frames. It is important to note that similar residual intensities can arise from different types of motion, which are likely to be associated with distinct distortion types.
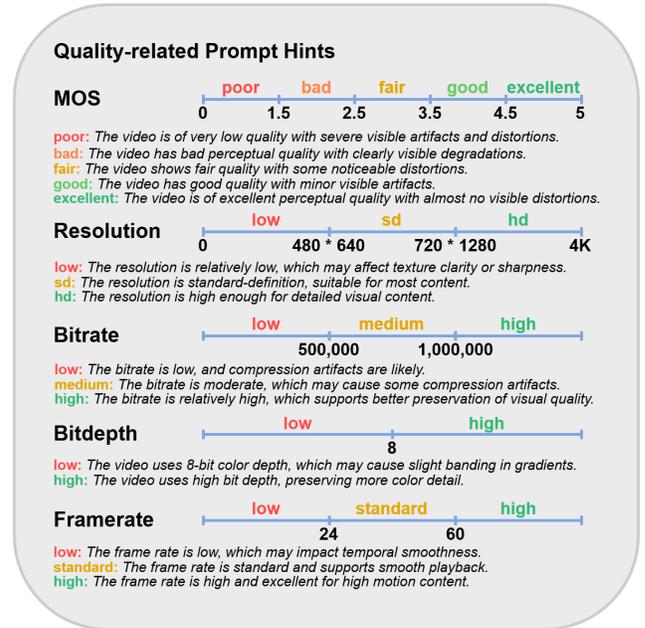
## B. Semantic-aligned Feature Extraction



Figure 2. Quality-related prompt hints derived from video metadata.

Based on video metadata across multiple dimensions, we derive a set of quality-related prompt hints to guide the input prompt, as illustrated in Fig. 2. The figure shows how metadata such as resolution, bitrate, and frame rate are mapped to qualitative levels based on their values, each

accompanied by descriptive hints. These mappings provide a structured way to integrate perceptual video quality factors into prompt construction for specific videos. Notably, we avoid disclosing ground-truth MOS scores. During training, quality-related hints are generated from quantized quality levels (ranging from poor to excellent) derived from metadata, guiding prompt text generation without including explicit scores in the descriptions. During inference, the quality level is set to prediction mode, eliminating the need for MOS.

Upon these hints, we design different quality-aware prompt settings, as illustrated in Fig. 3. The quality prompt is used to extract general perceptual captions at the frame level across multiple visual attributes. The fragment prompt focuses on localized degradations, and the residual prompt captures temporal differences between consecutive frames, jointly providing artifact captions at the fragment level. In the ablation study, we also include a content prompt to analyze the role of content descriptions in low-level video quality assessment. Each prompt is tailored to capture a different aspect of perceptual quality.

## C. Spatio-temporal Feature Extraction

### C.1. Spatial Vision Extractor

The Swin Transformer (SwinT) [3] employs a hierarchical sliding window mechanism and utilizes local self-attention to capture spatial features and long-range dependencies within images. To make full use of its strong spatial modeling capacity, we discard the classification head and retain only the backbone as our Spatial Vision Extractor (SVE) module.

Input frames of a video clip in the form of a tensor $\mathbf{X}_{\text{frame}} \in \mathbb{R}^{N \times C \times H \times W}$, where $N$ indicates the number of frames, $C$ is the number of channels, and $H$ and $W$ are the spatial dimensions. Each frame tensor is encoded by the SwinT backbone network $\phi_s$, and compact spatial representations are generated through GAP:

$$z_{\text{swint}} = \text{GAP}(\phi_{\text{swin}}(\mathbf{X}_{\text{frame}})) \in \mathbb{R}^{N \times d_s}, \quad (2)$$

where $d_s$ denotes the dimension of the spatial features. SVE module captures the spatial structural information of video frames, enabling spatial awareness for subsequent temporal modeling and multimodal fusion.

### C.2. Temporal Motion Extractor

To enhance the model's capacity for spatio-temporal feature modeling in videos, we designed a Temporal Motion Extractor (TME) module based on the SlowFast [2] architecture. TME uses a dual pathway to extract features across different temporal scales, enabling effective capture of quality fluctuations in videos while maintaining efficiency.

We construct two pathways with different temporal divisions by using an input frame sequence tensor $\mathbf{X}_{\text{sequence}} \in$ $\mathbb{R}^{B \times C \times T \times H \times W}$, where $B$ is batch size, $C$ is channels, $T$ is temporal length, and $H, W$ are the height and width. The slow pathway extracts low-frequency changes over longer temporal spans by downsampling at a rate of $r = \frac{1}{4}$, denoted as $X_{\text{slow}} = \text{Sample}(\mathbf{X}_{\text{sequence}}, r)$. The fast pathway retains the original frame rate, i.e., $X_{\text{fast}} = \mathbf{X}_{\text{sequence}}$, and serves to capture rapidly high-frequency details. The two pathways are fed into feature extractors $\phi_s$ and $\phi_f$, respectively. Spatio-temporal GAP is then applied to produce compact feature representations:

$$\mathbf{z}_{\text{slow}} = \text{GAP}(\phi_s(\mathbf{X}_{\text{slow}})), \ \mathbf{z}_{\text{fast}} = \text{GAP}(\phi_f(\mathbf{X}_{\text{fast}})). \quad (3)$$

Finally, the two feature vectors are concatenated to form the temporal feature:

$$\mathbf{z_{slowfast}} = [\mathbf{z}_{\text{slow}}; \mathbf{z}_{\text{fast}}] \in \mathbb{R}^{B \times (d_s + d_f)}, \quad (4)$$

where $d_s$ and $d_f$ are the feature dimensions of the slow and fast channels, respectively. This module extracts temporal motion features at different time scales, enhancing the model's sensitivity to temporal quality variations.

## D. Loss Functions

We adopted a composite loss function [4] to enhance accuracy and ranking consistency simultaneously. The first part of the loss function is the precision loss $\mathcal{L}_{\text{p}}$, which quantifies the mean absolute difference between the predicted regression score $\hat{y}\_i$ and the ground truth (MOS) $y_i$. It is defined as:

$$\mathcal{L}_{\text{p}} = \frac{1}{N} \sum_{i=1}^{N} |\hat{y}_i - y_i|. \quad (5)$$

The second is the ranking loss $\mathcal{L}_{\text{r}}$, which preserves ordinal consistency by comparing the relative rankings and computing the pairwise differences between sample pairs. It is defined as:

$$\mathcal{L}_{\text{r}} = \frac{1}{N^2} \sum_{i=1}^{N} \sum_{j=1}^{N} \max\left(0, |y_i - y_j| - s_{ij} \cdot (\hat{y}_i - \hat{y}_j)\right), \quad (6)$$

where $N$ is the number of videos, and $i$ and $j$ represent the indices of video samples. The sign weight $s_{ij}$ represents the ranking relationship inferred from the ground truth scores and is defined as $s_{ij} = 1$ if $y_i \geq y_j$, and $s_{ij} = -1$ if $y_i < y_j$. The composite loss $\mathcal{L}_c$ is formulated as a weighted combination of the two components: $\mathcal{L}_c = \lambda_1 \cdot \mathcal{L}_{\text{p}} + \lambda_2 \cdot \mathcal{L}_{\text{r}}$. Here, $\lambda \in \mathbb{R}$ balances the two loss terms: $\lambda_1$ improves prediction accuracy, while $\lambda_2$ enhances ordinal consistency. The loss is a weighted combination with $\lambda_1 = 0.6$ and $\lambda_2 = 1$.

**Quality Prompt (+ hint)**

**Quality Prompt:**

*You are a video quality assessment expert. Analyze this video frame and describe its perceptual quality in a short paragraph. Your analysis must address each of the following quality attributes:*

- *Sharpness (e.g., sharp, slightly fuzzy, very blurry)*
- *Focus (e.g., in-focus, out-of-focus)*
- *Noise (e.g., noiseless, mild noise, severe noise)*
- *Motion Blur (e.g., clear-motion, blur-motion)*
- *Flicker (e.g., stable, shaky)*
- *Compression Artifacts (e.g., blurring, ringing, blocking, banding)*
- *Color Issues (e.g., natural, faded)*
- *Exposure (e.g., well-exposed, overexposed, underexposed)*
- *Any other noticeable distortions (e.g., ghosting, flickering)*

*For each attribute, clearly indicate:*

- *Type of distortion (if any)*
- *Severity (choose from: none, mild, moderate, severe)*

*Respond with a short paragraph describing the perceptual quality based on the above, like: 'mild blur in background, moderate blocking in flat areas.' Do not describe the image content or name any objects and animals.*



**Fragment Prompt (+ hint)**

**Fragment Prompt:**

*This image is a small fragment cropped from a video frame and may lack full visual context. List the 1-2 most clearly visible quality issues (e.g., 'blurring', 'ghosting', 'flickering', 'sharpness', 'color inconsistency', or 'noise'). Only return distinct quality-related keywords. Do not describe the scene, do not guess objects, and do not repeat terms.*



**Residual Prompt (+ hint)**

**Residual Prompt:**

*This image is a residual fragment between the current and previous video frames, emphasizing visual differences caused by quality degradations. List the 1–2 most visible types of degradation (e.g., 'blurring', 'blockiness', 'ghosting', 'flickering', 'sharpness', 'color inconsistency', or 'noise'). Only return distinct quality-related keywords. Do not describe the scene, do not guess objects, and do not repeat terms.*

**Content Caption:** *bengal cat licking his paws*

**Content Prompt**

*Describe the visible content of this video frame as if explaining to someone who cannot see it. Mention the key objects, scene elements, or actions. Be concise and avoid any reference to image quality or technical terms.*

Figure 3. Different quality-aware prompt settings: quality prompt, fragment prompt, and residual prompt. We also include a content prompt for the ablation study.
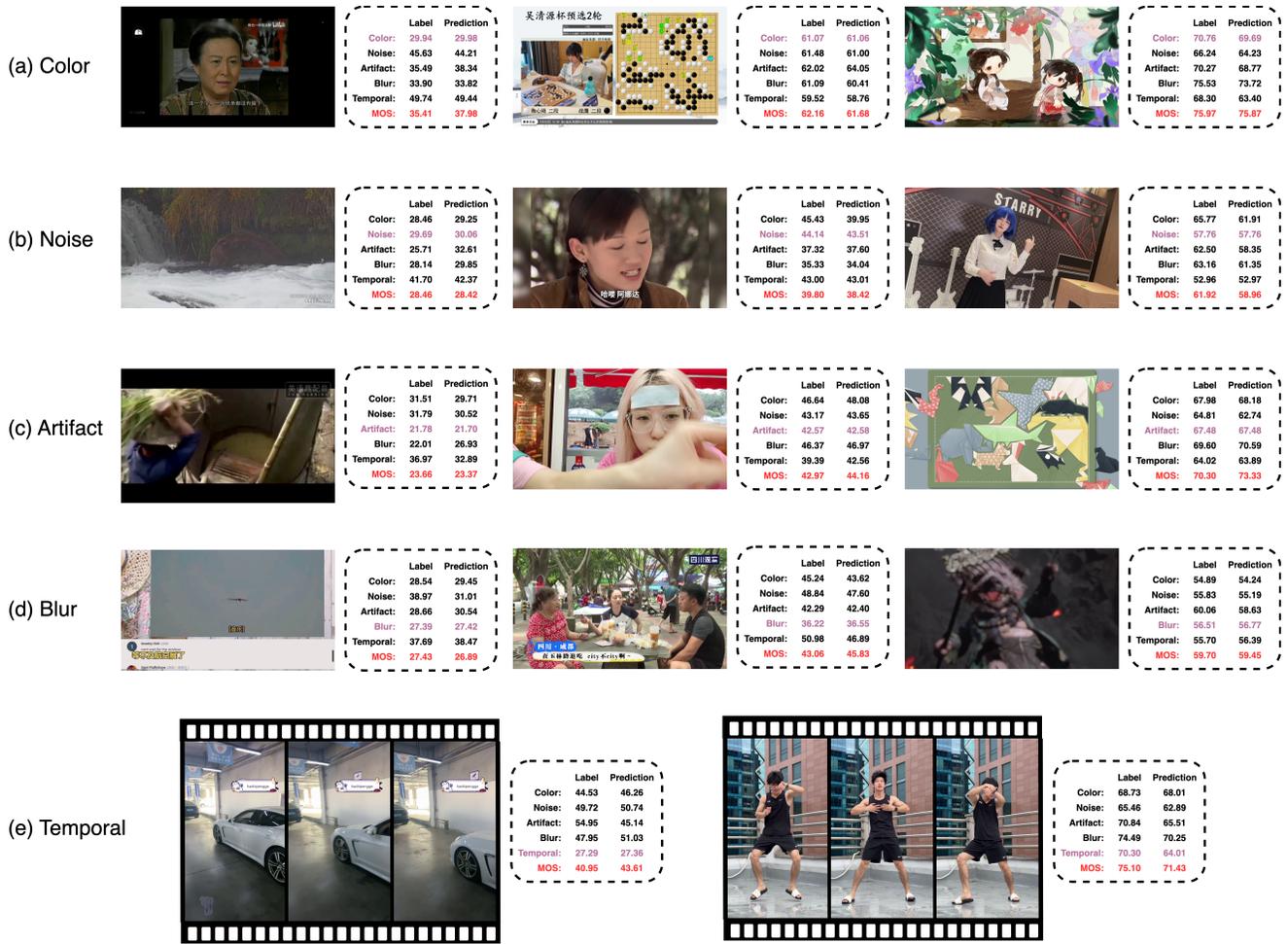
Figure 4. Predicted scores of CAMP-VQA and ground-truth MOS of FineVD [1] across different quality dimensions

# E. Quality Scoring On Different Dimensions

Building on our quality-aware captioning, we examined CAMP-VQA's impact on video quality scoring across different dimensions of the FineVD [1] dataset. Example videos with predicted CAMP-VQA scores and ground-truth FineVD MOS are shown in Figs. 4.

# References

[1] Huiyu Duan, Qiang Hu, Jiarui Wang, Liu Yang, Zitong Xu, Lu Liu, Xiongkuo Min, Chunlei Cai, Tianxiao Ye, Xiaoyun Zhang, et al. Finevq: Fine-grained user generated content video quality assessment. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 3206–3217, 2025. 4

[2] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6202–6211, 2019. 2

[3] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021. 2

[4] Shaoguo Wen and Junle Wang. A strong baseline for image and video quality assessment. *arXiv preprint:2111.07104*, 2021. 2