

A REAL-TIME BLIND QUALITY-OF-EXPERIENCE ASSESSMENT METRIC FOR HTTP ADAPTIVE STREAMING

Chunyi Li¹, May Lim², Abdelhak Bentaleb³, and Roger Zimmermann²

Shanghai Jiao Tong University¹, National University of Singapore², Concordia University³

ABSTRACT

In today’s Internet, HTTP Adaptive Streaming (HAS) is the mainstream standard for video streaming, which switches the bitrate of the video content based on an Adaptive BitRate (ABR) algorithm. An effective Quality of Experience (QoE) assessment metric can provide crucial feedback to an ABR algorithm. However, predicting such real-time QoE on the client side is challenging. The QoE prediction requires high consistency with the Human Visual System (HVS), low latency, and blind assessment, which are difficult to realize together. To address this challenge, we analyzed various characteristics of HAS systems and propose a non-uniform sampling metric to reduce time complexity. Furthermore, we design an effective QoE metric that integrates resolution and rebuffering time as the Quality of Service (QoS), as well as spatiotemporal output from a deep neural network and specific switching events as content information. These reward and penalty features are regressed into quality scores with a Support Vector Regression (SVR) model. Experimental results show that the accuracy of our metric outperforms the mainstream blind QoE metrics by 0.3, and its computing time is only 60% of the video playback, indicating that the proposed metric is capable of providing real-time guidance to ABR algorithms and improving the overall performance of HAS.

Index Terms— Quality of Experience, HTTP Adaptive Streaming, Blind Quality Assessment;

1. INTRODUCTION

Nowadays, video has become the dominant application on the Internet. According to Cisco’s survey [1], video services already consume more than 80% of current Internet capacity and demand is still growing. To meet the challenges posed by the transmission of large volumes of video data, content providers often use HTTP Adaptive Streaming (HAS), which can adapt to dynamic network conditions and various device resolutions. HAS delivers media content in small segments over the HTTP/TCP protocol stack, and because of its adoption by leading content providers, HAS has become the dominant delivery method for Video on Demand (VoD) services.

In a HAS client (*i.e.*, the media player), an Adaptive BitRate (ABR) algorithm selects a suitable bitrate level for the download of each video segment. The goal of ABR algorithms [2] is to maximize the user’s Quality of Experience

(QoE). Therefore, an effective objective QoE metric is crucial in HAS systems. It is either used after a video session ends to evaluate the performance of an ABR algorithm or during playback to guide the ABR algorithm in selecting the most appropriate bitrate demand for the next video segment.

Designing an effective objective QoE metric for bitrate guidance is challenging due to the following three requirements. *(i) High Consistency with HVS*: The QoE metric needs to be consistent with HVS so that it can accurately predict the user’s perceived quality of video playback. *(ii) Low Latency*: To provide guidance during transmission, the feedback of the QoE metric should be computed along with the video playback, and not after transmission like in ABR performance evaluation. In this case, the computation time of the QoE metric cannot be longer than the segment’s playback time, and hence its complexity should be reduced to ensure real-time feedback. Finally, *(iii) Blind Assessment*: The real-time prediction mentioned above is performed on the client. Hence, this should be a No-Reference (NR) task that uses only the compressed/distorted video available at the client, instead of Full-Reference (FR) or Reduced-Reference (RR) scenarios which use the original uncompressed source video or some of its reference features that have to be separately acquired from the server.

2. RELATED WORK AND CONTRIBUTIONS

Generally speaking, existing QoE metrics that provide guidance for HAS delivery can be classified into three types [3]: QoS-based, signal fidelity/content-based, and hybrid metrics.

QoS-based metrics [4, 5] have the lowest complexity as their computation typically involves a simple combination of network and/or client-side data, such as video bitrate, startup latency, and rebuffering duration. However, such metrics tend to show the least consistency with HVS as the video content is not considered in its computation, which has a strong influence on perceptual quality.

Content-based metrics [6, 7] analyze the video content and its signal fidelity using image or video quality assessment (I/VQA) metrics to predict a video’s distortion level [8]. IQA metrics take selected video frames as input and build a time series model [9] based on the quality of these frames to output a quality score for the video. Since an IQA metric needs to be computed repeatedly for each frame, it may not scale well. Existing VQA metrics also tend to have long run times. Only

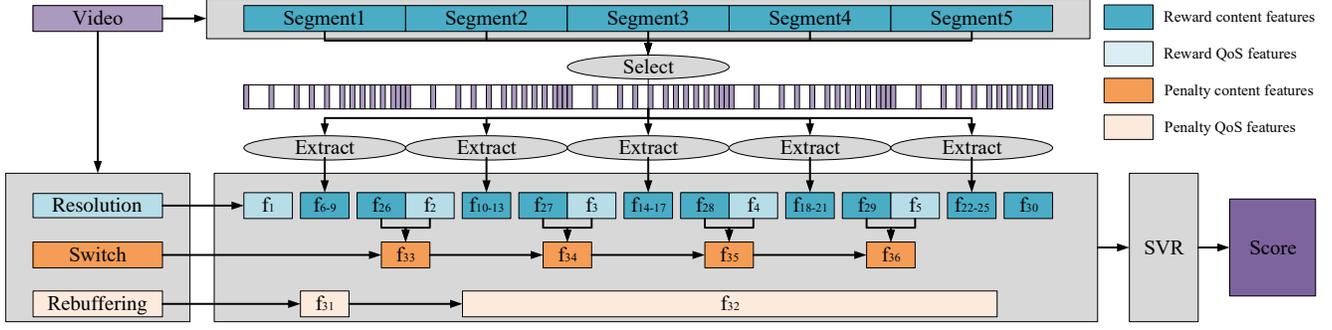


Fig. 1. The framework of the proposed method.

a few simplified NR-VQA metrics [10, 11] are able to provide real-time feedback for ABR, and such simplification leads to some degradation in their performance.

Hybrid metrics [12, 13] combine data from both the QoS and video content. Their QoE score is generally composed of a reward function represented by IQA/VQA metrics and a penalty function based on QoS factors. This allows the metric to reach a balance between time complexity and consistency with HVS, which is promising for its use in bitrate guidance. Unfortunately, some of the existing hybrid metrics have IQA/VQA kernels that are designed for FR tasks only [12], while the universal models [14] tend to show high consistency with HVS only on the FR or RR kernels but low consistency when it comes to NR. Therefore, the performance of such hybrid NR metrics can be further improved, as we show in this work.

As none of the existing metrics above can satisfy the three requirements of Section 1 altogether, we propose a new QoE assessment metric for HAS with the following contributions. (i) We perform a non-uniform sampling scheme since HVS has an increasing focus tendency throughout each segment. Without analyzing too many frames at the start of a segment, gradually increasing the sampling rate can reduce time complexity. (ii) We introduce QoS into reward features. As NR functions are not as effective as FR/RR, some QoS information can help our model perform better for blind assessment. (iii) We introduce content into penalty features. By analyzing specific frame changes instead of QoS fluctuations, the user's QoE can be better characterized.

3. PROPOSED METHOD

To design a QoE metric that can effectively meet the requirements discussed in Section 1, we identified novel features and adaptations that contribute materially to these requirements. The framework of our blind QoE metric is shown in Fig. 1 and includes three parts: *sampling*, *feature extraction*, and *quality regression*. Taking inspiration from hybrid metrics in Section 2, we extract four types of features: reward/penalty QoS features and reward/penalty content features. For content features, each video segment in the client's buffer is first non-uniformly sampled to select a suitable subset of image frames

for feature extraction. The QoS/content features are extracted by ResNet-50, texture analysis, and inter-frame difference, and then regressed through a support vector regression (SVR) model to give a quality score representing the user's current QoE. We discuss the details of each step below.

3.1. Sampling

For real-time QoE assessment, content-based/hybrid metrics analyze only a subset of frames to reduce complexity. While traditional metrics sample each segment uniformly, it is generally believed that frames within a segment may not contribute equally to QoE. To study their respective contributions, we run Brisque [6], a widely used NR-VQA metric, on the Waterloo sQoE III video dataset [15]. As each segment is two seconds long and an intra-coded frame usually appears every one second [16], we divide a segment into two halves (start/end) for non-uniform sampling and represent its QoE as: $QoE = w_s QoE_s + w_e QoE_e$, where w_s, w_e are weight parameters for the QoE of the start/end of a segment QoE_s and QoE_e , respectively. The Spearman Rank-order Correlation Coefficient (SRoCC) is used as the correlation function \mathcal{S} between QoE and the single-stimulus mean opinion scores (MOS) M obtained from the subjective assessment on the dataset: $\mathcal{S}(QoE, M) = w_s \mathcal{S}(QoE_s^{fr_s}, M) + w_e \mathcal{S}(QoE_e^{fr_e}, M)$, where fr_s and fr_e are the sampled frames from the start/end of a segment.

To better understand the relationship between sampling rate and correlation performance, we first used eight different sampling rates and three common IQA metrics (Niqe [17], Pique [18], Brisque [6]) to predict QoE. The normalized SRoCC between predicted QoEs and subjective scores show that the correlation factor is logarithmically related to the sampling rate as seen in Fig. 2. Hence, the sampling scheme can be transformed into the following optimization problem:

$$\begin{cases} fr = fr_s + fr_e \\ \text{maximize}(w_s \log(fr_s) + w_e \log(fr_e)) \end{cases} \quad (1)$$

where fr is the total number of sampled frames desired. Via Lagrange multiplier, the derivative of the log function in (1) gives fr as proportional to w for the best sampling scheme:

$$\frac{fr_s}{fr_e} = \frac{w_s}{w_e} = \frac{\mathcal{S}(\text{brisque}(seg_s), M)}{\mathcal{S}(\text{brisque}(seg_e), M)} \quad (2)$$

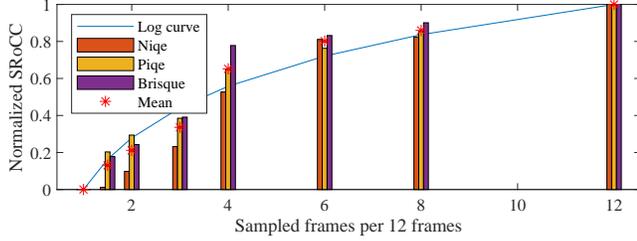


Fig. 2. Approximate convex logarithmic relationship between sampling rate and the normalized SRoCC.

where seg_s and seg_e are the start/end of a segment and QoE is predicted by $brisque(\cdot)$. The detailed proportion and derivation are attached in the supplementary.

3.2. Feature Extraction

We discuss how the four types of features are extracted below.

Reward QoS Feature. QoS refers to the basic network or client-side metrics during a video streaming session. Among them, video bitrate, quantization parameter (QP), frame rate, and resolution (height/width) have a positive impact on the user’s QoE. After studying the correlation performance of these factors, we found video height performs well as a reward feature r_1 to characterize the perceptual quality. Hence, $r_1 = \text{height}(seg)$, where seg is a video segment in HAS.

Reward Content Features. After sampling the frames in Section 3.1, we analyze this series of images which has three types of attainable features: structural, temporal, and chrominance. Structural features in images can reflect perceptual quality very well [19] and are widely used for QoE prediction. Temporal features can be computed independently or from the structural features via a spatial-temporal fusion. Chrominance features, like structural features, are computed independently from images. However, as the HVS processes visual signals in three channels, computing chrominance features may increase complexity considerably. To meet the real-time requirement, we include structural features and integrate temporal features with spatial-temporal fusion, while abandoning chrominance features.

The structural features can be divided into spatial and texture features. Spatial features refer to the relative spatial positioning or orientation of different elements in an image. An overview of our spatial feature extraction method is shown in Fig. 3. ResNet-50 is the backbone of the module, which can represent the spatial correlation between pixels and has proven to be quality-aware [20]. For the i -th sampled frame in a segment, we transform it into a gray map $g(i)$ as the input to ResNet-50, which extracts four features $L_i = l_{1\sim 4}(i)$ as:

$$\begin{cases} P(i) = \text{ResNet50}(g(i)) \\ [l_1(i), l_2(i)] = [\max(P(i)), \min(P(i))] \\ [l_3(i), l_4(i)] = [\text{avg}(P(i)), \text{std}(P(i))] \end{cases} \quad (3)$$

where $P(n)$ is the output from ResNet-50 and i is the frame index. On a segment level, instead of computing the four spa-

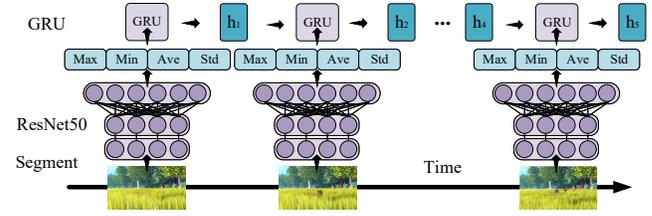


Fig. 3. The spatial feature extraction method.

tial features as a global average, we introduce a gated recurrent unit (GRU) [21] to capture the temporal relation between the spatial features. The spatiotemporal rewards $r_{2\sim 5}$ are:

$$\begin{cases} h_i = \text{GRU}(W \cdot L_i + b, h_{i-1}) \\ r_{2\sim 5} = h_{fr} \end{cases} \quad (4)$$

where W and b are the weights and bias parameters in the GRU [21], while h_i is a four-cell meta-array like L_i that restores the memory-forgetting mechanism of HVS in the time domain. Thus, these four features combine both spatial and temporal information.

Texture features refer to the surface characteristics of objects within images. Although operators such as Sobel, Laplace [22] are commonly used for texture extraction, given the complexity that ResNet-50 has already introduced for spatial features above, a simpler solution is needed here. Similar to how MacroBlocks (MBs) are used as the basic processing unit in video compression [23], we first compute the average row and column intensity values Ra_y, Ca_x from its gray map. Then we divide the gray map $g(i)$ into multiple 16×16 MBs. For each MB, we calculate the difference between the gray map value and the above-average values in the respective directions (horizontally, vertically, and diagonally). We select the minimum difference above as an MB’s texture, and combine them into the texture feature r_6 as:

$$\begin{cases} Hor_j = \sum |Ra_y(i) - g_{x,y}(i)| \\ Ver_j = \sum |Ca_x(i) - g_{x,y}(i)| \\ Dia_j = \sum |0.5(Ca_x(i) + Ra_y(i)) - g_{x,y}(i)| \\ r_6 = \sum \min(Hor_j, Ver_j, Dia_j) \end{cases} \quad (5)$$

where $Hor, Ver,$ and Dia are the texture information computed in three directions using the row average Ra_y and column average Ca_x , while j is the MB index.

Penalty QoS Features. Among QoS metrics, rebuffering has one of the largest negative impacts on QoE [24]. The rebuffering duration and number of rebuffering events are common penalty features used in QoE models [25, 12]. Besides, players tend to bear an initial buffering (of fewer than two seconds) in exchange for less rebuffering during playback [26]. We also note that the negative impact on QoE scales linearly with both rebuffering duration and number of rebuffering events and so we only consider one of them. Hence, we include the initial buffering duration and average rebuffering duration as penalty features p_1 and p_2 :

$$\begin{cases} p_1 = D_1 \\ p_2 = \frac{1}{T-1} \sum_{t=2}^T D_t \end{cases} \quad (6)$$

Table 1. The list of feature groups and their definitions.

Feature Group	Component	Origin	Description
Reward QoS	$f_{1\sim 5}$	r_1	Resolution
Reward Content	$f_{6\sim 25}$	$r_{2\sim 5}$	Spatial-Temporal
	$f_{26\sim 30}$	r_6	Texture
Penalty QoS	$f_{31\sim 32}$	$p_{1\sim 2}$	Rebuffering
Penalty Content	$f_{33\sim 36}$	p_3	Switching, Rebuffering (content-aware)

where T is the total number of segments, t is the segment index, and D_t is the buffering duration of the segment.

Penalty Content Features. Research into video QoS analysis [15] has shown that bitrate switching and rebuffering events can create a poor experience for users, the impact of which depends on several factors, including: (i) Switching pattern: It is generally believed [12] that dropping from high to low bitrate creates a more undesirable experience than going from low to high bitrate, and a long rebuffering duration [27] can further amplify the negative effects of such switching event. (ii) Video content: For a video where objects are moving slowly with minimal scene changes, the impact of a bitrate switching or rebuffering event is limited, while in an action movie, if the switching/rebuffering occurs during intense motion or when the scene has just changed, it generally leads to a significant drop in user’s QoE. The negative impact of each switching (or rebuffering) event is traditionally calculated using the change in inter-segment bitrate (or rebuffering duration), without considering the video content. Conversely, we use the efficient FastSSIM [28] metric to represent the inter-frame difference in video content between segments. The penalty content features p_3 combine all three factors above, namely the rebuffering time, switching level, and content mentioned in supplementary as shown in (7).

$$\begin{cases} swh = \text{ReLU}(\text{bitrate}(seg_t) - \text{bitrate}(seg_{t+1})) \\ p_3 = (1 + \frac{D_t}{C_1})(1 + \frac{swh}{C_2}) / \text{ssim}(seg_t, seg_{t+1}) \end{cases} \quad (7)$$

where swh is the inter-segment bitrate differential mapped to a ReLU function as it can characterize bitrate decline in the switching pattern mentioned above, while C_1 and C_2 are constants for normalization. The FastSSIM function $\text{ssim}(\cdot)$ is performed on the last sampled frame in segment seg_t and the first sampled frame in seg_{t+1} .

3.3. Overall QoE Regression

After the above operations, we finally obtain two reward feature groups r_1 and $r_{2\sim 6}$, and two penalty feature groups $p_{1,2}$ and p_3 to represent the positive and negative impact on user’s QoE, respectively. The feature groups are mapped to features $f_{1\sim 36}$ as shown in Table 1. Using these features extracted from QoS and video content information, a quality prediction model is constructed via support vector regression (SVR) to generate the overall QoE score as shown in Fig. 1. As suggested in prior studies [15], we use the most recent five segments in the playback for prediction in the SVR model. The SVR model is implemented using LIBSVM [29] with a radial basis function (RBF) kernel for feature fusion[30].

4. PERFORMANCE EVALUATION

4.1. Experiment Setup

The proposed metric is validated on the Waterloo sQoE III [15] and the LIVE Netflix II [35] datasets, which contain various subjectively-rated videos of diverse content types and video codecs, and streamed over various network conditions and ABR algorithms. The subjective study was done using dual-task single-stimulus (SS) experiments with user ratings provided as ground truth. The dataset is split randomly in an 80/20 ratio for training/testing while ensuring the same video content falls into the same set [10]. The partitioning and evaluation process is repeated 1,000 times for a fair comparison, and the average result is reported as the final performance.

We evaluate our metric in the following ways: First, to evaluate its consistency with HVS, we use three common correlation functions, namely SROCC, Kendall rank-order correlation coefficient (KROCC), and Pearson linear correlation coefficient (PLCC), to measure how well our metric correlates with the subjective scores. Second, to evaluate its latency, we measure the computation time of our metric on a laptop with an i7-8750H CPU, which is a common client specification for HAS services [1]. Third, since this is a blind assessment scenario, we compare our metric to 18 mainstream blind QoE assessment metrics as baselines. For model training, we left the parameters of ResNet-50 unchanged (which has been pre-trained for image classification on ImageNet [36]), and update the other parameters using the Adam optimizer [37]. The other learning-based methods are trained with similar mechanisms for a fair comparison. In the IQA models, we evaluated three different sampling rates by uniformly sampling every 5, 10, and 20 frames; in the hybrid models, as their original SSIM feature is not attainable for the NR task, we uniformly sample every 20 frames and apply the widely used Brisque [6] metric in their content-based features.

4.2. Experiment Results and Discussion

Table 2 and Fig. 4 show the performance results of the baseline and proposed methods, where a larger correlation factor indicates a higher consistency with the HVS. Computation time is calculated as a ratio of the video duration, where a smaller ratio indicates lower complexity and a value below 1 is needed to meet the real-time requirement. Among the content-based metrics, V-BLIINDS [31] has the best prediction performance but requires a computation time that is more than $40\times$ longer than the video itself, resulting in its inability to provide real-time feedback for bitrate guidance. The QoS-based metrics have much faster computation times of less than 1% of the video duration, but the predicted QoE has a relatively poor correlation with subjective scores. Most of the hybrid models achieve a better balance between consistency and latency. However, they are generally designed for FR features but their prediction performance becomes less ideal when

Table 2. Performance results on the Waterloo sQoE III and LIVE Netflix II datasets.

Type	Subtype	Method	Waterloo sQoE III				LIVE Netflix II			
			SROCC	KROCC	PLCC \uparrow	Time \downarrow	SROCC	KROCC	PLCC \uparrow	Time \downarrow
Content	IQA	Brisque5[6]	0.4959	0.3468	0.4690	1.944	0.3146	0.2182	0.2009	1.614
		Brisque10[6]	0.4842	0.3394	0.4600	0.965	0.2704	0.1880	0.1344	1.276
		Brisque20[6]	0.4547	0.3146	0.4478	0.498	0.2543	0.1794	0.1517	0.455
		Niqe10[17]	0.4021	0.2732	0.4488	1.027	0.6538	0.4761	0.6684	0.529
		Piqe10[18]	0.4232	0.2807	0.4349	1.128	0.6746	0.4898	0.6871	0.541
	VQA	VIIDEO[7]	0.3946	0.2651	0.4903	8.047	0.2843	0.1885	0.3228	5.624
		V-BLIINDS[31]	0.7389	0.5456	0.7244	45.625	0.7510	0.5755	0.7653	42.755
		Resnet50[32]	0.5707	0.4113	0.5635	0.381	0.4278	0.2995	0.4317	0.307
	FAST-VQA[11]	0.7391	0.5500	0.7710	0.246	0.5137	0.3644	0.5748	0.232	
QoS	Client	FTW[4]	0.1835	0.1337	0.3229	0.001	0.0804	0.0858	0.0648	0.001
		MoK2011[25]	0.1687	0.1294	0.2156	0.001	0.0795	0.0650	0.0874	0.001
	Network	Liu2012[33]	0.2529	0.1717	0.2424	0.001	0.6633	0.4684	0.6366	0.001
		Xue2014[34]	0.3412	0.2245	0.3081	0.003	0.5830	0.4123	0.4961	0.003
		Yin2015[5]	0.1458	0.0932	0.3232	0.007	0.0804	0.0616	0.0648	0.007
Hybrid	Mix	SQI[13]	0.1515	0.1100	0.2225	0.501	0.7347	0.5298	0.6329	0.458
		TV-QoE[14]	0.5068	0.3565	0.4667	0.524	0.6686	0.4136	0.5109	0.482
		Bentaleb2016[12]	0.1979	0.1387	0.3405	0.498	0.4454	0.2982	0.4530	0.456
		KSQI[3]	0.5285	0.3875	0.5268	0.505	0.7394	0.5492	0.7315	0.462
		Proposed	0.8627	0.6871	0.8824	0.606	0.7739	0.5914	0.7898	0.691

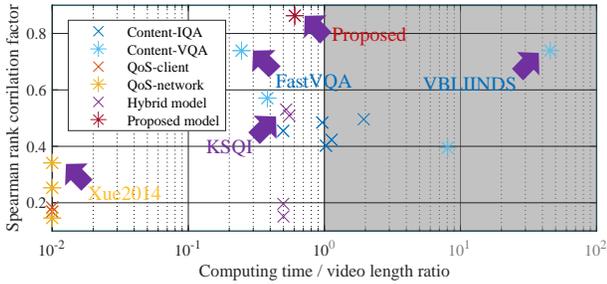


Fig. 4. SROCC and time ratio performance of the baseline and proposed methods. Metrics with time ratio < 1 (white panel) can realize real-time QoE prediction for bitrate guidance.

switched to NR. Results show that our metric outperforms all QoS-based and hybrid metrics in all three correlation measures for both datasets (with a gain of up to 0.35 against other hybrid metrics), and outperforms all content-based metrics in SROCC and PLCC for both datasets while keeping the computation time ratio below 1. In terms of global QoE metric performance factors[38], our Area Under the ROC Curve (AUC) value is about 0.09 ahead of SOTA HAS QoE metric and outperforms all current metrics for correct classification. Hence, our metric can overcome the absence of reference information in blind assessment scenarios and still provide effective QoE predictions that have high consistency with HVS (with an average SROCC of 0.81) and sufficiently low latency (of about 65% of the video playback duration).

4.3. Ablation Study

To validate the contributions of our sampling method and the different feature types, we also conduct an ablation study and its results are shown in Table 3. The factors are specified as: (1) Non-uniform sampling, (2) Reward QoS features, (3) Re-

Table 3. Performance results of abandoning different features on the Waterloo sQoE III dataset.

Abandoned	SROCC	KROCC	PLCC	Time
None	0.8627	0.6871	0.8824	0.606
(1)	0.8514	0.6716	0.8694	0.605
(2)	0.4521	0.3270	0.4761	0.606
(3) $r_2 \sim r_5$	0.7662	0.5822	0.8194	0.381
(3) r_6	0.8468	0.6680	0.8637	0.441
(3) All	0.7313	0.5478	0.7490	0.216
(4)	0.8319	0.6544	0.8691	0.606
(5)	0.8434	0.6767	0.8639	0.390
(2)(4)	0.3204	0.2214	0.4006	0.605
(3)(5)	0.7010	0.5215	0.7073	0.001

ward content features (both spatial and textural), (4) Penalty QoS features, and (5) Penalty content features. The results show that removing any single factor leads to performance degradation, which confirms that they all contribute to the performance results in Table 2. The time cost for each group can ensure a real-time assessment. Ablation results also show that QoS feature groups (2)(4) are more effective than the content feature groups (3)(5), and a combination of them achieves desirable performance, while the sampling method (1) further enhances the model’s performance with little extra time cost.

5. CONCLUSIONS

In this study, we target the challenge of measuring perceptual quality in HAS clients for bitrate guidance. A blind QoE assessment metric is proposed to provide QoE feedback of high consistency with HVS, at low latency, and without the use of reference information. Specifically, we map QoS and video content information to both reward and penalty features and include a non-uniform sampling mechanism to identify relevant frames for analysis. Experiments show that the proposed

metric achieves the best results in two databases across three correlation measures, which suggests strong consistency with HVS, while meeting latency and blind assessment requirements. This metric is suited to perform real-time QoE assessment on the client side, which can help improve the overall resource utilization of HAS services.

6. REFERENCES

- [1] Cisco, “Cisco Visual Networking Index: Forecast and Trends, 2018–2023,” *White Paper*, 2020.
- [2] A. Bentaleb, B. Taani, A. C. Begen, C. Timmerer, and R. Zimmermann, “A Survey on Bitrate Adaptation Schemes for Streaming Media Over HTTP,” *IEEE Communications Surveys & Tutorials*.
- [3] Z. Duanmu, W. Liu, D. Chen, Z. Li, Z. Wang, et al., “A knowledge-driven quality-of-experience model for adaptive streaming videos,” *arXiv preprint arXiv:1911.07944*, 2019.
- [4] T. Hoßfeld, R. Schatz, E. Biersack, and L. Plissonneau, “Internet video delivery in YouTube: From traffic measurements to quality of experience,” in *Data Traffic Monitoring and Analysis*. Springer, 2013.
- [5] X. Yin, A. Jindal, V. Sekar, and B. Sinopoli, “A control-theoretic approach for dynamic adaptive video streaming over HTTP,” in *ACM SIGCOMM*, 2015.
- [6] A. Mittal, A. K. Moorthy, and A. C. Bovik, “No-reference image quality assessment in the spatial domain,” *IEEE TIP*, 2012.
- [7] A. Mittal, M. A. Saad, and A. C. Bovik, “A completely blind video integrity oracle,” *IEEE TIP*, 2015.
- [8] G. Zhai and X. Min, “Perceptual image quality assessment: a survey,” *Science China Information Sciences*, 2020.
- [9] X. Sui, K. Ma, Y. Yao, and Y. Fang, “Perceptual quality assessment of omnidirectional images as moving camera videos,” *IEEE TVCG*, 2021.
- [10] Z. Tu, X. Yu, Y. Wang, N. Birkbeck, B. Adsumilli, and A. C. Bovik, “RAPIQUE: Rapid and accurate video quality prediction of user generated content,” *IEEE OJSP*, 2021.
- [11] H. Wu, C. Chen, J. Hou, L. Liao, A. Wang, W. Sun, Q. Yan, and W. Lin, “FAST-VQA: Efficient end-to-end video quality assessment with fragment sampling,” *ECCV*, 2022.
- [12] A. Bentaleb, A. C. Begen, and R. Zimmermann, “SDNDASH: Improving QoE of HTTP adaptive streaming using software defined networking,” in *ACM Multimedia*, 2016.
- [13] Z. Duanmu, K. Zeng, K. Ma, A. Rehman, and Z. Wang, “A quality-of-experience index for streaming video,” *IEEE Journal of Selected Topics in Signal Processing*, 2016.
- [14] D. Ghadiyaram, J. Pan, and A. C. Bovik, “Learning a continuous-time streaming video QoE model,” *IEEE TIP*, 2018.
- [15] Z. Duanmu, A. Rehman, and Z. Wang, “A quality-of-experience database for adaptive video streaming,” *IEEE Transactions on Broadcasting*, 2018.
- [16] H. Amirpour, E. Çetinkaya, C. Timmerer, and M. Ghanbari, “Fast multi-rate encoding for adaptive http streaming,” in *Data Compression Conference (DCC)*, 2020.
- [17] A. Mittal, R. Soundararajan, and A. C. Bovik, “Making a “completely blind” image quality analyzer,” *IEEE Signal Processing Letters*, 2012.
- [18] N. Venkatanath, D. Praneeth, M. C. Bh, S. S. Channappayya, and S. S. Medasani, “Blind image quality evaluation using perception based features,” in *IEEE NCC*, 2015.
- [19] Q. Li, W. Lin, J. Xu, and Y. Fang, “Blind image quality assessment using statistical structural and luminance features,” *IEEE TMM*, 2016.
- [20] J. Kim, H. Zeng, D. Ghadiyaram, S. Lee, L. Zhang, and A. C. Bovik, “Deep convolutional neural models for picture-quality prediction: Challenges and solutions to data-driven image quality assessment,” *IEEE Signal Process Mag*, 2017.
- [21] J. Yan, J. Li, Y. Fang, Z. Che, X. Xia, and Y. Liu, “Subjective and objective quality of experience of free viewpoint videos,” *IEEE TIP*, 2022.
- [22] M. Sonka, V. Hlavac, and R. Boyle, *Image processing, analysis, and machine vision*, Cengage Learning, 2014.
- [23] X. Min, J. Zhou, G. Zhai, P. Le Callet, X. Yang, and X. Guan, “A metric for light field reconstruction, compression, and display quality evaluation,” *IEEE TIP*, 2020.
- [24] C. G. Bampis and A. C. Bovik, “Learning to predict streaming video qoe: Distortions, rebuffering and memory,” *arXiv preprint arXiv:1703.00633*, 2017.
- [25] R. K. Mok, X. Luo, E. W. Chan, and R. K. Chang, “QDASH: a QoE-aware DASH system,” in *ACM MMSys*, 2012.
- [26] Tisa-Selma, A. Bentaleb, and S. Harous, “Video qoe inference with machine learning,” in *IEEE IWCMC*, 2021.
- [27] B. Taraghi, M. Nguyen, H. Amirpour, and C. Timmerer, “Intense: In-depth studies on stall events and quality switches and their impact on the quality of experience in http adaptive streaming,” *IEEE Access*, 2021.
- [28] M.-J. Chen and A. C. Bovik, “Fast structural similarity index algorithm,” *Journal of Real-Time Image Processing*, 2011.
- [29] C.-C. Chang and C.-J. Lin, “Libsvm: a library for support vector machines,” *ACM TIST*, 2011.
- [30] L. Krasula, K. Fliegel, and P. Le Callet, “Fftmi: Features fusion for natural tone-mapped images quality evaluation,” *IEEE TMM*, 2020.
- [31] M. A. Saad, A. C. Bovik, and C. Charrier, “Blind prediction of natural video quality,” *IEEE TIP*, 2014.
- [32] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *IEEE CVPR*, 2016.
- [33] X. Liu, F. Dobrian, H. Milner, J. Jiang, V. Sekar, I. Stoica, and H. Zhang, “A case for a coordinated internet video control plane,” in *ACM SIGCOMM*, 2012.
- [34] J. Xue, D.-Q. Zhang, H. Yu, and C. W. Chen, “Assessing quality of experience for adaptive HTTP video streaming,” in *IEEE ICME Workshops*, 2014.
- [35] C. G. Bampis, Z. Li, I. Katsavounidis, T.-Y. Huang, C. Ekanadham, and A. C. Bovik, “Towards perceptually optimized end-to-end adaptive video streaming,” *arXiv preprint arXiv:1808.03898*, 2018.
- [36] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al., “Imagenet large scale visual recognition challenge,” *IJCV*, 2015.
- [37] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [38] L. Krasula, K. Fliegel, P. Le Callet, and M. Klíma, “On the accuracy of objective image and video quality models: New methodology for performance evaluation,” in *QoMEX*, 2016.