

Cut-FUNQUE: An Objective Quality Model for Compressed Tone-Mapped High Dynamic Range Videos

Abhinav K. Venkataramanan, Cosmin Stejerean, Ioannis Katsavounidis,
Hassene Tmar, and Alan C. Bovik, *Life Fellow, IEEE*

Abstract—High Dynamic Range (HDR) videos have enjoyed a surge in popularity in recent years due to their ability to represent a wider range of contrast and color than Standard Dynamic Range (SDR) videos. Although HDR video capture has seen increasing popularity because of recent flagship mobile phones such as Apple iPhones, Google Pixels, and Samsung Galaxy phones, a broad swath of consumers still utilize legacy SDR displays that are unable to display HDR videos. As result, HDR videos must be processed, i.e., tone-mapped, before streaming to a large section of SDR-capable video consumers. However, server-side tone-mapping involves automating decisions regarding the choices of tone-mapping operators (TMOs) and their parameters to yield high-fidelity outputs. Moreover, these choices must be balanced against the effects of lossy compression, which is ubiquitous in streaming scenarios. In this work, we develop a novel, efficient model of objective video quality named Cut-FUNQUE that is able to accurately predict the visual quality of tone-mapped and compressed HDR videos. Finally, we evaluate Cut-FUNQUE on a large-scale crowdsourced database of such videos and show that it achieves state-of-the-art accuracy.

Index Terms—High Dynamic Range, Tone Mapping, Video Quality, Perceptual Uniformity

I. INTRODUCTION

The Human Visual System (HVS) encounters a diverse range of luminances, or brightness, in real-world scenarios, spanning from the faint glow of starlight at 0.0003 cd/m^2 (nits) to the intense brightness of sunlight reaching up to 30,000 nits on a clear day. Thanks to the adjustment of the pupil size by the iris to varying brightness levels, the HVS can perceive an extensive range of luminances, ranging from approximately 10^{-6} nits to 10^8 nits. However, conventional imaging and display systems are limited to capturing or producing narrow ranges of luminances, typically up to about 100 nits. These systems, categorized as low or standard dynamic range (SDR) systems, are also only capable of capturing or displaying about 35% of the visible color gamut. Notable examples of legacy SDR standards include sRGB [1] and ITU BT. 709 [2].

Over the years, the development of high dynamic range (HDR) imaging has aimed to better align imaging and display systems with the capabilities of the HVS. Contemporary HDR standards, exemplified by ITU BT. 2100 [3], have the capacity to capture luminances spanning from 10^{-4} to 10^4 nits, along

with a wide color gamut (WCG) that encompasses approximately 75% of the visible color volume. To achieve this, HDR imaging employs computational imaging techniques that blend two or more images taken at different exposure settings. Novel “optoelectrical transfer functions” (OETFs) extend the conventional notion of “gamma” from legacy Cathode Ray Tube (CRT) displays to effectively encode and transmit this wide range of brightnesses.

In particular, the BT. 2100 standard incorporates two encoding functions, namely the Perceptual Quantizer (PQ) [4] and the Hybrid Log-Gamma (HLG) [5]. PQ is engineered as a “forward-compatible” standard, capable of encoding luminances up to 10^4 nits, and is commonly utilized by professional studios to deliver high-quality HDR content. Notably, the PQ encoding function serves as the foundation of standards like HDR10 [6] and HDR10+ [7]. In contrast, HLG is designed to be “backward-compatible” with SDR standards by incorporating a “gamma” curve similar to that used in SDR, in the SDR luminance range. While the HLG standard does not explicitly define a peak luminance, a nominal value of 1000 nits is often adopted. Due to its backward compatibility, HLG has gained adoption by satellite TV networks for the delivery of HDR content [8]. Notably, the emerging Dolby Vision standard [9] supports both PQ and HLG encoding functions.

The widespread streaming of HDR videos faces a significant challenge due to the limited availability of true HDR displays. According to the BT. 2100 standard, true HDR systems are defined as those capable of supporting a minimum of 1000 nits [3]. However, a majority of budget-friendly HDR displays fall short of this criterion, reaching a peak luminance of 800 nits or less [10]. In fact, a substantial portion of existing displays are only capable of supporting SDR formats. Consequently, to ensure accessibility of HDR video content to a broader consumer base, it becomes essential to “down-convert” them to the SDR range, a process commonly referred to as “tone-mapping.”

Section IV-A provides a brief overview of various tone-mapping methods proposed in the literature. However, the inherent limitations of SDR systems lead to distortions during the tone-mapping process. These distortions typically manifest as either contrast reduction or enhancement, as well as the potential losses or amplification of details in dark or bright regions. Furthermore, the remapping of color across dynamic ranges and the wide color gamut (WCG) employed by high dynamic range (HDR) can give rise to chromatic

This research was sponsored by a grant from Meta Video Infrastructure, and by grant number 2019844 from the National Science Foundation AI Institute for Foundations of Machine Learning (IFML).

distortions, including hue shifts and chroma-clipping artifacts [11] [12]. Additionally, the necessity for lossy compression when streaming videos over the internet introduces another layer of distortions, including flattening of details, blocking in regions containing high motion or texture and banding on smoother regions.

Here, we tackle the problem of objective quality assessment of compressed tone-mapped videos. In particular, we design the new Cut-FUNQUE objective quality model which is based on the recently developed FUNQUE [13] [14] framework, and demonstrate its effectiveness in automatically predicting the subjective quality of tone-mapped videos. Cut-FUNQUE is the sum total of three key novel contributions.

- 1) The first contribution is the development of a novel perceptually uniform encoding of color signals (**PUColor**) that we use to represent both HDR and SDR color stimuli in a common domain. In this manner, PUColor enables the meaningful comparison of stimuli across dynamic ranges, which is essential when comparing HDR and SDR videos.
- 2) Secondly, Cut-FUNQUE utilizes a **binned-weighting** approach to separately handle image regions having different visual characteristics such as brightness, contrast, and temporal complexity.
- 3) Finally, Cut-FUNQUE also utilizes novel **statistical similarity** measures of visual quality to overcome the limitations of pixel-wise comparisons across dynamic ranges.

We demonstrate the efficacy of Cut-FUNQUE by conducting evaluations on the recently introduced LIVE Tone-Mapped HDR (LIVE-TMHDR) subjective database.

The remainder of this manuscript is organized as follows. Section II provides an overview of the literature regarding the objective quality assessment of tone-mapped images and videos. Section III describes the proposed Cut-FUNQUE model in detail, and in Section IV, we report the results of comparing Cut-FUNQUE against existing video quality models in the literature.

Finally, we present a summary of our findings and directions for future work in Section V.

II. BACKGROUND

Due to the longstanding nature of the HDR tone-mapping problem, several algorithms have been developed to assess the quality of tone-mapped media, particularly images. The Dynamic Range Independent Quality Assessment (DRIQA) model [15] was the first attempt at comparing HDR and SDR images despite their different dynamic ranges. The model generates distortion maps using contrast sensitivity models, cortex transforms, and psychometric probability models to measure loss, gain, and reversal of visible contrast. However, by design, the algorithm does not provide accurate “single number” predictions of visual quality from the distortions maps (such as by averaging) [15].

The Tone Mapped Quality Index (TMQI) [16] was the first quality model designed to conduct tone-mapped image quality prediction. It introduced the popular framework of hybrid

Full Reference (FR) - No Reference (NR) methods. TMQI combines measures of structural fidelity, similar to Structural Similarity (SSIM) [17], and naturalness, which may be considered an NR model, to predict overall quality. A variant of TMQI, called TMQI-NSS [18], incorporates visual attention into the structure term and uses Natural Scene Statistics (NSS) to quantify naturalness. The NSS model borrows from NR algorithms such as BRISQUE [19] and NIQE [20].

The FSITM [21] algorithm is a modified version of the FSIM [22] quality model, and it is based on analyzing the local phase of complex log-Gabor wavelet transform coefficients. Binarized phase values from the reference HDR and test SDR images are compared to obtain a structural similarity measure, while naturalness is defined as the same similarity metric applied between the SDR image and the log of the HDR image. Combining this model with TMQI has been shown to improve performance [21]. FFTMI [23] is a fusion-based model built on this principle, and its “atom quality models” include FSITM computed from RGB channels, TMQI, and a no-reference model of naturalness [24]. These atom features are combined to predict a single quality score using a regressor.

Finally, the Tone Mapped Video Quality Index (TMVQI) [25] is an adaptation of the TMQI model. TMVQI uses an updated contrast sensitivity model to modulate contrast estimates, and an updated naturalness model that was estimated from HDR video frames.

Among no-reference models, the Blind TMQI (BTMQI) [26] model utilizes entropy estimates from simulated multi-exposure images to predict structural fidelity and borrows its naturalness measure from TMQI. HIGRADE [27] employs a Natural Scene Statistics (NSS) approach similar to algorithms such as BRISQUE [19] by modeling the statistics of Mean Subtracted Contrast Normalized coefficients. In addition, HIGRADE incorporates structural information by modeling the statistics of gradient structure tensors.

The RcNet [28] model is a deep-learning NR quality model that uses a suite of Convolutional Neural Networks to predict the DRIQA distortions maps without using a reference HDR source. The distributions of MSCN coefficients of the SDR image and the predicted distortion maps are used to predict its quality. The state-of-the-art deep method is the Multi-Scale Multi-Layer [29] model, which uses pooled features extracted from the intermediate layers of a pretrained ResNet-50 network. The outputs of three hidden layers are concatenated to yield a feature vector containing 9216 features, and partial least squares (PLS) is used to reduce the final dimensionality of the feature vector to fifteen.

III. CUT-FUNQUE

As described in Section II, successful models of tone-mapped video quality are hybrid ones that include both full-reference and no-reference features, such as the structure and naturalness terms in TMQI [16]. Moreover, with the increase in available dataset sizes, feature-based models have emerged as a powerful tool for quality modeling, both for no-reference models such as V-BLIINDS [30], VIDEVAL [31], and HIGRADE [27], and full-reference methods such as VMAF

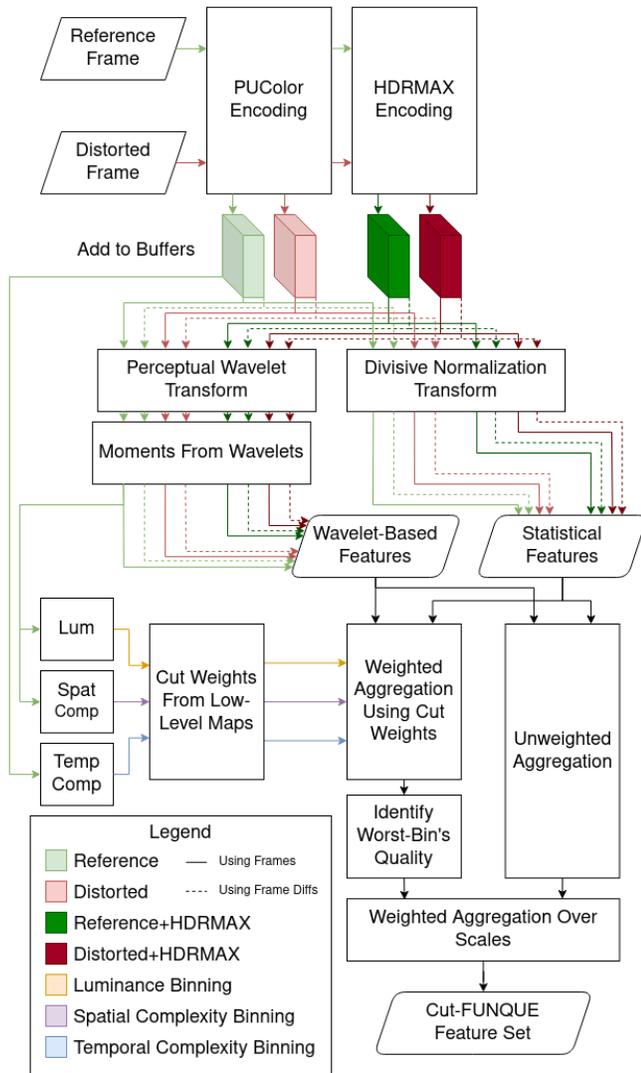


Fig. 1. An overview of the Cut-FUNQUE quality prediction model.

[32] and the FUNQUE framework [14] [33]. Therefore, we employ a fusion-based bag-of-features approach consisting of smaller “atom” quality models and features designed based on perceptual principles.

To design atom features appropriate for the quality assessment of tone-mapped videos, we first made a set of observations regarding the challenges faced by prior work in the area, and the effects of tone-mapping and compression. The primary challenge that is endemic to the field of tone-mapped video quality assessment is identifying a common domain in which both HDR and SDR may be well represented. Based on prior work in the area of perceptually uniform representations [34], we posited that a good transform domain designed for this purpose would be **perceptually uniform over a wide range of brightnesses**.

While PU21 [34] is a promising choice for such a transform, it was derived based on a luminance contrast sensitivity model. Since color distortions are an important aspect of tone mapping, an ideal transform would be a perceptually uniform color encoding function (PUColor) that can represent

three-channel color values. We describe our proposed PUColor encoding function in Section III-A.

Furthermore, we make the following observations regarding the nature of distortions that arise from tone-mapping and compression.

- Since tone-mapping involves applying compressive or sigmoidal non-linearities, **contrast changes at various luminance levels may be different**. These non-linearities are typically responsible for distortions, such as loss of detail/contrast, in very bright and dark regions of a video. On the other hand, while the contrast in mid regions is typically preserved, distortions in the mid regions may be indicative of global brightness changes.
- **Loss/gain of detail in various regions may be affected by both tone-mapping and compression**. Compression alters local contrast by inducing blockiness (which is blur-like) in textured (high-contrast) regions and banding in smooth (low-contrast) regions. On the other hand, local processing in tone-mapping often boosts regions of low local contrast and attenuates regions of high local contrast. Therefore, the nature of distortions also varies with the amount of local contrast.
- **Regions with high temporal variation are also affected by both tone-mapping and compression**, in the form of flickering and misalignment/blocking due to the quantization of motion vectors.

To incorporate this prior knowledge into the algorithm, we partitioned each video frame into patches (referred to as “cuts” since they are non-overlapping) and (softly) classified each cut into various “bins” depending on their low-level spatio-temporal properties. We then employed a weighted aggregation method to pool quality features within each bin to characterize different “types” of frame regions. Finally, we identified the (predicted) worst-quality regions to account for quality-based saliency. This procedure is described in further detail in Section III-F. Furthermore, we combined the proposed PUColor with the recently developed HDRMAX transform [35], which emphasizes quality prediction on bright and dark regions of input frames, which are perceptually important for HDR stimuli.

To enable effective quality prediction, we extracted various quality-aware features from input videos. These features included spatial and temporal models of quality such as SSIM [17], VIF [14], DLM [36], and ST-RRED [37]. Notably, all the aforementioned quality models are “pixel-based” methods that directly compare (processed) pixel values. Such methods are limited in their capacity to accurately compare video frames across dynamic ranges.

For example, the SSIM score between a pair of image regions is unity only if their pixel values are identical. However, since pixel values of an SDR frame occupy a limited range, they are likely to be significantly different in value from their HDR frame counterparts, irrespective of the quality of tone-mapping. To overcome this, we propose novel statistical similarity (StSIM) measures of quality, based on principles of Natural Scene Statistics (NSS). StSIM features compare the distributions of normalized pixel values, rather than comparing

pixel values directly. This makes StSIM features more robust to changes in dynamic range.

The feature set employed by Cut-FUNQUE is described in further detail in Sections III-D and III-E. A flowchart illustrating the Cut-FUNQUE processing flow is presented in Figure 1.

A. A Perceptually Uniform Encoding Function for Color

In this section, we present the derivation of a perceptually uniform color encoding function (PUColor) using a “post-receptoral” model of the spatio-chromatic contrast sensitivity function (SCCSF) [38]. This SCCSF model has been used in prior work to derive perceptually uniform luminance encoding functions [34]. Here, we design the PUColor encoding function to be approximately uniform along three “chromatic” directions, which we choose to correspond to opponent color channels - achromatic, red-green, and blue-yellow.

If t denotes the detection threshold, i.e., a Just-Noticeable-Difference (JND), we say that the *PUColor* encoding function of LMS values $\mathbf{x} = [x_L, x_M, x_S]^T$ is uniform along three “chromatic” vectors \mathbf{u}_i if

$$PUColor(\mathbf{x} + t(\mathbf{x}, \mathbf{u}_i) \cdot \mathbf{u}_i) - PUColor(\mathbf{x}) \approx C\mathbf{e}_i, \quad (1)$$

where C is an arbitrary constant and the set $\{\mathbf{e}_i\}$ refers to the standard basis vectors. Here, “LMS” refers to a color space that represents the responses of cone cells in the retina in the long, medium, and short-wavelength regions. In the analysis below, we use $C = 1$ for convenience. We then scale the derived function such that the range of encoded values is in the range $[0, 1]$.

Using a first-order approximation of *PUColor*, the uniformity condition in Eq. 1 may be rewritten as

$$\mathbf{J}(\mathbf{x})\mathbf{u}_i \approx \frac{1}{t(\mathbf{x}, \mathbf{u}_i)}\mathbf{e}_i, \quad (2)$$

where $\mathbf{J}(\mathbf{x})$ denotes the Jacobian matrix of the partial derivatives of *PUColor*. It may be observed that Eq. 2 is a generalized eigenvalue relationship between the threshold function and the direction of uniformity. Collecting the chromatic vectors into a matrix $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \mathbf{u}_3]$, then

$$\mathbf{J}(\mathbf{x}) = \text{diag}(t^{-1}(\mathbf{x}, \mathbf{u}_i))\mathbf{U}^{-1}, \quad (3)$$

where *diag* denotes a diagonal matrix.

To complete the description of the Jacobian matrix, the threshold function t must be characterized. Here, we utilized the threshold function proposed in the SCCSF model [38], which was defined in terms of a transformation matrix \mathbf{M}_{ARB} , three “base sensitivity functions” $s_A(\cdot), s_R(\cdot), s_B(\cdot)$, and spatial frequency ρ

$$t(\mathbf{x}, \mathbf{u}) = \min_{\rho} \left\| \left(\frac{\mathbf{s}_{\text{ARB}}(\rho, x_L + x_M)}{x_L + x_M} \right) \odot (\mathbf{M}_{\text{ARB}}\mathbf{u}) \right\|_2^{-1}, \quad (4)$$

where \odot denotes the elementwise (Hadamard) product.

Finally, to obtain the *PUColor* encoding function, we enforced the boundary condition that $PUColor(\mathbf{0}) = 0$ and

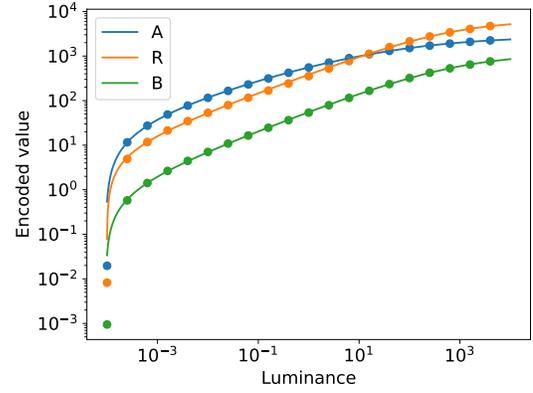


Fig. 2. Numerical integrals of the threshold function and approximations using non-linear functions

integrated the Jacobian over the straight line $z = \{\lambda\mathbf{x}; \lambda \in [0, 1]\}$. Thence,

$$PUColor(\mathbf{x}) = \int_0^1 \mathbf{J}(\mathbf{z}(\lambda))\mathbf{z}'(\lambda) d\lambda \quad (5)$$

$$= \left(\int_0^1 \frac{d\lambda}{t(\lambda\mathbf{x}, \mathbf{u}_i)} \right) \odot (\mathbf{U}^{-1}\mathbf{x}). \quad (6)$$

The three integrals, corresponding to the three chromatic directions \mathbf{u}_i , were evaluated numerically. A closer look at Eq. 4 reveals that t depends on \mathbf{x} only through $y = x_L + x_M$, which denotes the luminance. Consequently, the integrals in Equation 6 share the same property. Therefore, to apply the *PUColor* encoding in practice, we approximate the results of the numerical integrals using non-linear functions of the form $h(y)/y$, where

$$h(y; p) = \left(\frac{p_1 + p_2 y^{p_4}}{1 + p_3 y^{p_4}} \right)^{p_5}. \quad (7)$$

A comparison of the values of numerical integrals and the corresponding non-linear function fits for various luminance values is presented in Figure 2. All three numerical fits achieved r^2 values over 0.999. To ensure that the three “chromatic directions” correspond to an opponent representation of color, we chose $\mathbf{U} = \mathbf{M}_{\text{DKL}}^{-1}$, where \mathbf{M}_{DKL} is the transformation matrix corresponding to the Derrington-Krauskopf-Lennie (DKL) opponent color space [39]. Therefore, the final *PUColor* encoding function has the simple form

$$PUColor(\mathbf{x}) = \frac{\mathbf{h}_{\text{ARB}}(x_L + x_M)}{x_L + x_M} \odot (\mathbf{M}_{\text{DKL}} \cdot \mathbf{x}). \quad (8)$$

The encoded values consist of one achromatic channel (L) and two chromatic channels that may be denoted by a and b , analogous to color spaces such as CIELAB [40]. The two chromatic channels encode color together, where the “angle” between the two components $\arctan(b/a)$ encodes the hue, and the magnitude $\sqrt{a^2 + b^2}$ encodes chromaticity. Therefore, we combine the two channels into one complex-valued color channel $c = a + jb$, where $j = \sqrt{-1}$, and the argument ($\arg c$)

and magnitude ($|c|$) of c represent the hue and chromaticity respectively.

B. HDRMAX Encoding

HDRMAX [35] refers to a non-linear pre-processing method that seeks to emphasize perceptually important image regions in HDR videos. In particular, HDRMAX emphasizes the measurement of distortions in bright and dark image regions. This is achieved by pre-processing image frames in two steps - local normalization and a point non-linear transform.

Local min-max normalization is used to first normalize local luma (and chroma) values within overlapping 17×17 neighborhoods to the range $[-1, 1]$:

$$I_{minmax}(i, j) = 2 \frac{I(i, j) - I_{min}(i, j)}{I_{max}(i, j) - I_{min}(i, j)}. \quad (9)$$

Normalized luma (and chroma) values are then subjected to a double-exponential non-linearity

$$\text{HDRMAX}(x) = \text{sgn}(x) \left(e^{4|x|} - 1 \right). \quad (10)$$

Although HDRMAX was designed to enhance the prediction of the quality of HDR stimuli [33], it has also been effective for SDR video quality prediction [35]. So, we adopt it as a pre-processing step for both HDR reference and SDR (i.e., tone-mapped) test videos.

C. Perceptually-Sensitized Wavelet Transform

Following the FUNQUE framework [13] [14], we compute a perceptually-sensitized wavelet decomposition of all analyzed reference and test frames after applying the PUColor and HDRMAX transformations. The Self Adaptive Scale Transform (SAST) [41] was first applied to account for viewing distance, which was assumed to be three times the display height as in [14]. This assumption corresponds to typical HD TV viewing distances.

Then, a Haar wavelet transform was applied to obtain a band-pass decomposition of the input frames. To account for the visibility of various artifacts due to the contrast sensitivity of the visual system, the subbands of the Haar decomposition were weighted using Watson's model of wavelet subband contrast sensitivity [42]. We preferred Watson's model of contrast sensitivity over other models described in [14] due to its wavelet domain definition and its modeling of chroma visibility.

The perceptually-sensitized wavelet decomposition was used in [14] to compute multi-scale local moments (i.e., mean, variance, and covariance) of the reference and test frames, which were used to compute Multi-Scale Enhanced SSIM (MS-ESSIM) [43]. We adopt a similar procedure here, but we reuse the same moments to compute three sets of quality features - ESSIM, Visual Information Fidelity (VIF) [44], and Spatio-Temporal Reduced Reference Entropic Differences (ST-RRED) [37]. Moreover, local moments of the reference frame were also used to compute cut-assignment weights, which are described in Section III-F. Such extensive computation sharing follows the FUNQUE framework and contributes

to Cut-FUNQUE's efficiency despite the large number of features.

Let x and y denote the reference and test images, and $X_{l,\lambda,\theta}$ and $Y_{l,\lambda,\theta}$ denote the wavelet decompositions of their luma channels, where $\lambda \in \{1 \dots L\}$ denotes the wavelet level (i.e., "scale") and $\theta \in \{A, H, V, D\}$ denote the "orientation" of the subband. Then, local moments of the luma channel within non-overlapping windows of size $2^\lambda \times 2^\lambda$ are computed as

$$\mu_{l,x,\lambda}(i, j) = 2^{-\lambda} X_{l,\lambda,A}(i, j), \quad (11)$$

$$\mu_{l,y,\lambda}(i, j) = 2^{-\lambda} Y_{l,\lambda,A}(i, j), \quad (12)$$

$$\sigma_{l,x,\lambda}^2(i, j) = 2^{-2\lambda} \sum_{k=1}^{\lambda} \sum_{P_{ij}^k} \sum_{\{H,V,D\}} X_{l,k,\theta}^2(m, n), \quad (13)$$

$$\sigma_{l,y,\lambda}^2(i, j) = 2^{-2\lambda} \sum_{k=1}^{\lambda} \sum_{P_{ij}^k} \sum_{\{H,V,D\}} Y_{l,k,\theta}^2(m, n), \quad (14)$$

$$\sigma_{l,xy,\lambda}(i, j) = 2^{-2\lambda} \sum_{k=1}^{\lambda} \sum_{P_{ij}^k} \sum_{\{H,V,D\}} X_{l,k,\theta}(m, n) Y_{l,k,\theta}(m, n), \quad (15)$$

where $P_{ij}^k = \{(m, n) \mid i2^{\lambda-k} \leq m < (i+1)2^{\lambda-k}, j2^{\lambda-k} \leq n < (j+1)2^{\lambda-k}\}$ denotes disjoint $2^{\lambda-k} \times 2^{\lambda-k}$ blocks.

Analyzing Eq. (13), it may be seen that multi-scale variances can be computed iteratively using the relationship

$$\sigma_{l,x,\lambda+1}^2(i, j) = 2^{-2} \sum_{m=i}^{i+1} \sum_{n=j}^{j+1} \sigma_{l,x,\lambda}(m, n)^2 + \tilde{\sigma}_{l,x,\lambda+1}^2(i, j), \quad (16)$$

where

$$\tilde{\sigma}_{l,x,\lambda+1}^2(i, j) = 2^{-2(\lambda+1)} \sum_{\{H,V,D\}} X_{l,\lambda+1,\theta}^2(i, j). \quad (17)$$

The energies $\sigma_{l,x,\lambda+1}^2(i, j)$, and $\sigma_{l,xy,\lambda+1}(i, j)$ were also computed iteratively in a similar manner.

Local moments of the complex-valued chroma channel were computed similarly, with the key difference being the use of complex conjugates when computing variance and covariance. As an example,

$$\sigma_{c,xy,\lambda}(i, j) = 2^{-2\lambda} \sum_{k=1}^{\lambda} \sum_{P_{ij}^k} \sum_{\{H,V,D\}} X_{c,k,\theta}(m, n) Y_{c,k,\theta}^*(m, n), \quad (18)$$

where $*$ denotes the complex conjugate. Once again, multi-scale moments were computed iteratively as in Eq. (16).

Using this method, we compute local moments within 8×8 , 16×16 , 32×32 , and 64×64 windows, which correspond to the four scales at which Cut-FUNQUE is applied.

D. Wavelet-Domain Feature Extraction

In this section, we describe the wavelet-domain features computed using the perceptually-sensitized wavelet transform described in Section III-C. In particular, this section describes the procedure to calculate local quality-aware features that were spatially aggregated using the method described in Section III-F.

We computed four sets of wavelet-domain features - SSIM, VIF, RRED, and DLM, from which SSIM, VIF, and RRED were computed using the local moments described in Section III-C. Local SSIM scores were split into their two components, namely luminance and contrast-structure similarity [17]. Since we compute SSIM features on both luma and complex-valued chroma channels, we rename the two components as $SSIM_\mu$ and $SSIM_\sigma$ as general terminology we use for both channels. SSIM features from the luma channel were computed as:

$$\text{L-SSIM}_\mu(i, j) = \frac{2\mu_{l,x}(i, j)\mu_{l,y}(i, j) + C_1}{\mu_{l,x}^2(i, j) + \mu_{l,y}^2(i, j) + C_1} \quad (19)$$

$$\text{L-SSIM}_\sigma(i, j) = \frac{2\sigma_{l,xy}(i, j) + C_2}{\sigma_{l,x}^2(i, j) + \sigma_{l,y}^2(i, j) + C_2} \quad (20)$$

and the chroma SSIM features were computed as

$$\text{C-SSIM}_\mu(i, j) = \frac{2|\mu_{c,x}(i, j)\mu_{c,y}(i, j)| + C_1}{|\mu_{c,x}^2(i, j)| + |\mu_{c,y}^2(i, j)| + C_1} \quad (21)$$

$$\text{C-SSIM}_\sigma(i, j) = \frac{2|\sigma_{c,xy}(i, j)| + C_2}{\sigma_{c,x}^2(i, j) + \sigma_{c,y}^2(i, j) + C_2}. \quad (22)$$

This process was repeated across four scales, both with and without HDRMAX, yielding a total of 32 feature maps.

Visual Information Fidelity (VIF) [44] is defined as the ratio of two estimates of mutual-information corresponding to the test and reference video frames. As described in [44], VIF assumes that distortions arise from a channel modeled as

$$Y_{l,\lambda}(i, j) = g_{l,\lambda}(i, j)X_{l,\lambda}(i, j) + N_{l,\lambda}(i, j), \quad (23)$$

where $g_{l,\lambda}(i, j)$ is the gain of the distortion channel and $N_{l,\lambda}(i, j) \sim \mathcal{N}(0, \sigma_v^2)$ is additive white Gaussian noise.

We adopted a similar distortion model for the complex-valued chroma channel, with the modification that all quantities may take complex values and the noise channel is a complex Gaussian, i.e., $N_{c,\lambda}(i, j) \sim \mathcal{CN}(0, \sigma_v^2)$.

The local mutual information (MI) estimates, computed under a Gaussian assumption, were obtained from the luma channel's moments using the following equations:

$$g_{l,\lambda}(i, j) = \sigma_{l,xy,\lambda}(i, j) / \sigma_{l,x,\lambda}^2(i, j), \quad (24)$$

$$\sigma_{l,v,\lambda}^2(i, j) = \sigma_{l,y,\lambda}^2(i, j) - g_{l,\lambda,\theta}(i, j)\sigma_{l,xy,\lambda}(i, j), \quad (25)$$

$$\text{MI-Test}_{l,\lambda}(i, j) = \log \left(1 + \frac{g_{l,\lambda}^2(i, j)\sigma_{l,x,\lambda}^2(i, j)}{\sigma_{l,v,\lambda}^2(i, j) + \sigma_n^2} \right), \quad (26)$$

and

$$\text{MI-Ref}_{l,\lambda} = \log \left(1 + \frac{\sigma_{l,x,\lambda}^2(i, j)}{\sigma_n^2} \right), \quad (27)$$

where σ_n^2 denotes the noise variance of the assumed distortion channel. Hence, local VIF scores were computed as

$$\text{VIF}_{l,\lambda}(i, j) = \frac{\text{MI-Test}(i, j)}{\text{MI-Ref}(i, j)} \quad (28)$$

A similar procedure was used to obtain local VIF scores for the chroma channel under a complex Gaussian assumption. We omit the analogous expressions of the quantities described above for brevity, noting that differential entropy of a 1-D complex Gaussian distribution is identical to that of a 2-D real-valued Gaussian distribution having its real and imaginary parts as its components. Furthermore, we repeated the process on the differences of successive frames to obtain estimates of temporal quality, which we denote by "TVIF". Hence, we obtained four quality features L-VIF, C-VIF, L-TVIF, and C-TVIF, each repeated at four scales and with and without HDRMAX, yielding a total of 32 features.

The Reduced Reference Entropic Difference (RRED) quality features [37] are defined as the difference between scaled local entropy estimates, also obtained under a Gaussian assumption using the same local moments as SSIM and VIF:

$$h_{l,\lambda}(i, j) = \alpha_{l,\lambda}(i, j) \log (2\pi e(\sigma_{l,\lambda}^2(i, j) + \sigma_n^2)), \quad (29)$$

where the weighting factors are given by

$$\alpha_{\lambda,\theta}(i, j) = \log (1 + \sigma_{\lambda,\theta}^2(i, j)). \quad (30)$$

The difference map, i.e. spatial RRED (SRRED) quality map, is obtained as

$$\text{L-SRRED}(i, j) = |h_{l,x,\lambda}(i, j) - h_{l,y,\lambda}(i, j)| \quad (31)$$

Similarly, differences between local entropies of frame differences yield the temporal RRED (TRRED) quality map, and chroma SRRED and TRRED (C-SRRED and C-TRRED) maps were obtained by computing local entropies of the chroma channel under a complex Gaussian assumption. This procedure was repeated at four scales and for HDRMAX-transformed frames to yield a total of 32 RRED feature maps.

The final wavelet-domain feature corresponds to the Detail Loss Metric (DLM) [36]. A brief description of DLM is as follows. The first step in computing DLM involves applying a "decoupling step." The decoupling step is based on the following distortion model of the wavelet subband coefficients $\theta \in \{H, V, D\}$:

$$Y_{\lambda,\theta}(i, j) = \gamma_{\lambda,\theta}(i, j)X_{\lambda,\theta}(i, j) + A_{\lambda,\theta}(i, j), \quad (32)$$

where the gain factor γ models attenuation of local gradients due to detail loss, $R_{\lambda,\theta}(i, j) = \gamma_{\lambda,\theta}(i, j)X_{\lambda,1}(i, j)$ are the "restored" coefficients, and $A_{\lambda,\theta}(i, j)$ are the "additive impairments."

The "restored" wavelet decomposition are computed from the given frames via

$$\hat{R}_{\lambda,\theta}(i, j) = \hat{\gamma}_{\lambda,\theta}(i, j)X_{\lambda,\theta}(i, j), \quad (33)$$

where

$$\psi_{x,\lambda}(i, j) = \arctan \left(\frac{X_{\lambda,V}(i, j)}{X_{\lambda,H}(i, j)} \right), \quad (34)$$

$$\psi_{y,\lambda}(i, j) = \arctan \left(\frac{Y_{\lambda,V}(i, j)}{Y_{\lambda,H}(i, j)} \right), \quad (35)$$

$$\Delta\psi_\lambda(i, j) = |\psi_{x,\lambda}(i, j) - \psi_{y,\lambda}(i, j)|, \quad (36)$$

and

$$\hat{\gamma}_{\lambda,\theta}(i, j) = \begin{cases} \frac{Y_{\lambda,\theta}(i, j)}{X_{\lambda,\theta}(i, j)}, & \Delta\psi_\lambda(i, j) < 1^\circ \\ \text{clip}\left(\frac{Y_{\lambda,\theta}(i, j)}{X_{\lambda,\theta}(i, j)}, 0, 1\right), & \text{else} \end{cases}. \quad (37)$$

The quantity $\Delta\psi$ is used to preserve contrast enhancement, which scales both the horizontal and vertical subband coefficients of the test frame [36]. The additive impairment coefficients are then computed as:

$$\hat{A}_{\lambda,\theta}(i, j) = Y_{\lambda,\theta}(i, j) - \hat{R}_{\lambda,\theta}(i, j). \quad (38)$$

The additive impairments are used to mask the restored coefficients using the model

$$\tilde{R}_{\lambda,\theta}(i, j) = \left(\hat{R}_{\lambda,\theta}(i, j) - M_\lambda(i, j)\right)^+, \quad (39)$$

where

$$M_\lambda(i, j) = \sum_{\theta} \sum_{i-1}^{i+1} \sum_{j-1}^{j+1} w_{ij}(m, n) \left| \hat{A}_{\lambda,\theta}(m, n) \right|, \quad (40)$$

$$w_{ij}(m, n) = \frac{1 + \delta(m - i, n - j)}{30}, \quad (41)$$

$(\cdot)^+$ denotes clipping negative values to zero, and δ is the Kronecker delta function. Finally, modifying the definition in [36] to yield local scores, local DLM scores are computed as the following ratio:

$$\text{DLM}(i, j) = \frac{\left(\sum_{\theta \in \{H, V, D\}} \tilde{R}_{\lambda,\theta}(i, j)^3\right)^{1/3}}{\left(\sum_{\theta \in \{H, V, D\}} \tilde{X}_{\lambda,\theta}(i, j)^3\right)^{1/3}}. \quad (42)$$

E. Natural Scene Statistics Features

In addition to wavelet-domain features, we incorporated principles from natural scene statistics to compare input frames by their bandpass distributions. We achieved this by applying a bandpass transform called Mean Subtracted Contrast Normalization (MSCN), which has been used in a wide variety of popular NR quality prediction models, such as NIQE [20], BRISQUE [19], and HIGRADE [27].

The MSCN transform is a Divisive Normalization Transform (DNT) applied to either the mean-subtracted luma or chroma channel

$$\tilde{I}(i, j) = \frac{I(i, j) - \mu(i, j)}{\sigma(i, j) + \epsilon}, \quad (43)$$

where $\mu(i, j)$ and $\sigma(i, j)$ are the local and means and standard deviations. Furthermore, to model contrast naturalness, we also applied the MSCN transform to the $\sigma(i, j)$ field, yielding “ σ -MSCN” coefficients.

MSCN coefficients obtained in this manner have been observed to follow regular quality-aware statistical properties, which are captured by modeling their distribution using Generalized Gaussian Distributions (GGDs) [19] [27]. The GGD

centered at 0 may be defined using its probability density function (PDF) as:

$$f_{GGD}(x; \alpha, b) = \frac{\alpha}{2b\Gamma(\frac{1}{\alpha})} \exp\left(-\left(\frac{|x|}{b}\right)^\alpha\right), \quad (44)$$

where α is called the “shape parameter”, b is called the “scale parameter”, and Γ denotes the Gamma function, which is defined as follows:

$$\Gamma(x) = \int_0^\infty t^{x-1} e^{-t} dt. \quad (45)$$

We fit GGDs to MSCN and σ -MSCN coefficients by estimating distribution parameters using a moment-matching procedure [45], then use them as “global” no-reference NSS features. In addition, we also capture the second-order statistics of MSCN coefficients by modeling the distribution of products of spatially adjacent MSCN coefficients. We chose the neighboring pixel locations along four directions - horizontal, vertical, and two diagonals (H, V, D_1, D_2). As in [19] and [20], Asymmetric GGDs (AGGDs) are fit to the distributions of these products using a moment-matching approach [46] and their parameters used as quality-aware features. An AGGD centered at zero may be described by a PDF of the form

$$f_{AGGD}(x; \alpha, b_l, b_r) = \begin{cases} \frac{\alpha}{(b_l + b_r)\Gamma(\frac{1}{\alpha})} \exp\left(-\left(\frac{|x|}{b_l}\right)^\alpha\right) & x < 0, \\ \frac{\alpha}{(b_l + b_r)\Gamma(\frac{1}{\alpha})} \exp\left(-\left(\frac{|x|}{b_r}\right)^\alpha\right) & x \geq 0, \end{cases} \quad (46)$$

where α is a “shape” parameter and b_l, b_r are left and right “scale parameters.”

To limit the number of features, we average the quality-aware AGGD parameter estimates over the four directions to obtain a final set of second-order global NR NSS features.

The use of statistical fits to MSCN coefficients has typically been limited to the domain of NR quality assessment. However, the availability of a reference video opens up the possibility of defining “statistical similarity” (StSim) quality models, analogous to structural similarity models such as SSIM. Here, we introduce the novel method of comparing KL divergences between NSS distributions as a measure of quality-aware statistical dissimilarity.

Let a pair of corresponding GGD parameters obtained from a reference (HDR) and distorted (SDR) frame pair be (α_1, b_1) and (α_2, b_2) . Then, the first-order statistical dissimilarity (FOSD) is the KL Divergence between the two GGDs:

$$\text{FOSD} = \frac{\Gamma\left(\frac{\alpha_2 + 1}{\alpha_1}\right)}{\Gamma\left(\frac{1}{\alpha_1}\right)} \left(\frac{b_1}{b_2}\right)^{\alpha_2} - \frac{1}{\alpha_1} + \log\left(\frac{\alpha_1 b_2 \Gamma\left(\frac{1}{\alpha_2}\right)}{\alpha_2 b_1 \Gamma\left(\frac{1}{\alpha_1}\right)}\right). \quad (47)$$

Similarly, the second-order statistical dissimilarity (SOSD) between two frames described by AGGDs $(\alpha_1, b_{l1}, b_{r1})$ and

$(\alpha_2, b_{l2}, b_{r2})$ is:

$$\text{SOSD} = \frac{\Gamma\left(\frac{\alpha_2+1}{\alpha_1}\right)}{\Gamma\left(\frac{1}{\alpha_1}\right) (b_{l1} + b_{r1})} \left(\frac{b_{l1}^{\alpha_2+1}}{b_{l2}^{\alpha_2}} + \frac{b_{r1}^{\alpha_2+1}}{b_{r2}^{\alpha_2}} \right) - \frac{1}{\alpha_1} + \log \left(\frac{\alpha_1 (b_{l2} + b_{r2}) \Gamma\left(\frac{1}{\alpha_2}\right)}{\alpha_2 (b_{l1} + b_{r1}) \Gamma\left(\frac{1}{\alpha_1}\right)} \right). \quad (48)$$

We repeated the same procedure locally to obtain first and second-order NSS fits to coefficients within cuts of size 8×8 , 16×16 , 32×32 , and 64×64 , corresponding to analyzing the video at four scales. Then, FOSD and SOSD estimates were obtained for each cut, yielding local statistical dissimilarity estimates.

F. Binned Weighting to Isolate Frame Regions

As explained earlier, we expect different regions of a frame that vary in their characteristics to be affected differently by tone-mapping and compression distortions. So, we first partition the frame into ‘‘cuts’’, i.e., non-overlapping patches of equal sizes to enable local analysis. CutQA utilizes a multi-scale approach, with the finest scale corresponding to cuts of size 8×8 and the coarsest scale corresponding to 64×64 cuts.

Next, in accordance with the three observations made earlier, we characterize the **luminance** (brightness), **spatial complexity**, and **temporal complexity** of each cut using three low-level measures. The brightness of the cut is characterized by the mean luma value over the cut. Spatial complexity, i.e., contrast, is often characterized by the standard deviation of luma values within the cut. However, brighter patches tend to have higher standard deviations since their pixel values are greater. To account for this phenomenon, we instead used the coefficient-of-variation (CoV), defined as the ratio of standard deviation to mean. Finally, to characterize the temporal complexity, we computed the averaged temporal coefficient of variation of luma values over the previous four frames.

Formally, consider a cut c and an input HDR frame I_{HDR} , the three low-level measures are computed as

$$L(c) = E_{(i,j) \in c} [I_{HDR}(i, j)], \quad (49)$$

$$S(c) = \frac{Std_{(i,j) \in c} [I_{HDR}(i, j)]}{E_{(i,j) \in c} [I_{HDR}(i, j)]}, \quad (50)$$

and

$$T(c) = \frac{E_{(i,j) \in c} [Std_{t \in [-3,0]} [I_{HDR}(i, j)]]}{E_{(i,j) \in c} [E_{t \in [-3,0]} [I_{HDR}(i, j)]]}, \quad (51)$$

where $E_{(i,j) \in c}$ and $Std_{(i,j) \in c}$ denote computing means and standard deviations over spatial locations in a cut c , and $E_{t \in [-3,0]}$ and $Std_{t \in [-3,0]}$ denote computing mean and standard deviation over the four most recent frames, including the current frame. Local means and standard deviations used to compute $L(c)$ and $S(c)$ are computed from the perceptually-sensitized Haar wavelet decomposition described in Section III-C. The four most recent frames used to compute $T(c)$

are stored in a ‘‘reference frame buffer,’’ which is updated framewise.

The set of low-level descriptors obtained in this manner are then assigned to four equally spaced bins of each type, respectively called luminance, spatial, and temporal contrast (L-, S-, and T-) bins. In particular, for each bin type, a Gaussian membership function is used to define four membership weights for each cut. For example, consider a cut c with mean luminance $L(c)$, and let the four luminance bins have centers \tilde{L}_b and widths $w^{(L)}$ ($b = 0, 1, 2, 3$). The membership weights for the cut corresponding to each luminance bin are computed as

$$M_b^{(L)}(c) = \exp \left(-\frac{(L(c) - \tilde{L}_b)^2}{2(w^{(L)}/2)^2} \right). \quad (52)$$

This yields a soft classification of the cut into four luminance bins. Using the same procedure, each cut is also soft classified into four S-bins and four T-bins using weights $M_b^{(S)}(c)$ and $M_b^{(T)}(c)$ respectively, which are calculated as in Eq. (52).

An example of partitioning a video frame into 8×8 cuts and corresponding bin classification weights is presented in Fig. 3. The scene (Fig. 3a) consists of cars moving on a road surrounded by various buildings and lights. Figs. 3b-3e visualize L-bin membership weights $M_0^{(L)} - M_3^{(L)}$, Figs. 3f-3i visualize S-bin membership weights $M_0^{(S)} - M_3^{(S)}$, and Figs. 3j-3m visualize T-bin membership weights $M_0^{(T)} - M_3^{(T)}$.

In each weight map, a brighter yellow color denotes a higher membership weight, i.e., the corresponding region ‘‘belongs more’’ to the bin. From that figure, it may be observed that the luminance bins distinguish between bright and dark regions, while spatial complexity bins distinguish between detailed regions, and temporal complexity bins can segment moving objects such as cars.

Bin-membership weights are then used to aggregate local quality scores to bin-level quality scores. Bin-level quality scores represent the quality of different ‘‘types of regions’’ in the image. For example, quality scores aggregated within luminance bin 0 encode the quality of dark image regions. These scores are obtained using weighted averaging, which we illustrate using the example of SRRED. Let $SRRED(c)$ denote the SRRED feature value computed for each cut c . The weighted-aggregated SRRED values for the four luminance bins ($b = 0, 1, 2, 3$) are computed as

$$SRRED_b^{(L)} = \frac{\sum_c M_b^{(L)}(c) \times SRRED(c)}{\sum_c M_b^{(L)}(c)}. \quad (53)$$

Similarly, $SRRED_b^{(S)}$ and $SRRED_b^{(T)}$ are computed, and the process is repeated for all features.

The one exception is SSIM, which is aggregated using weighted Minkowski pooling, based on the recommendations in [43]. We preferred Minkowski pooling over the CoV pooling used in [13] because we found the ratio used in CoV

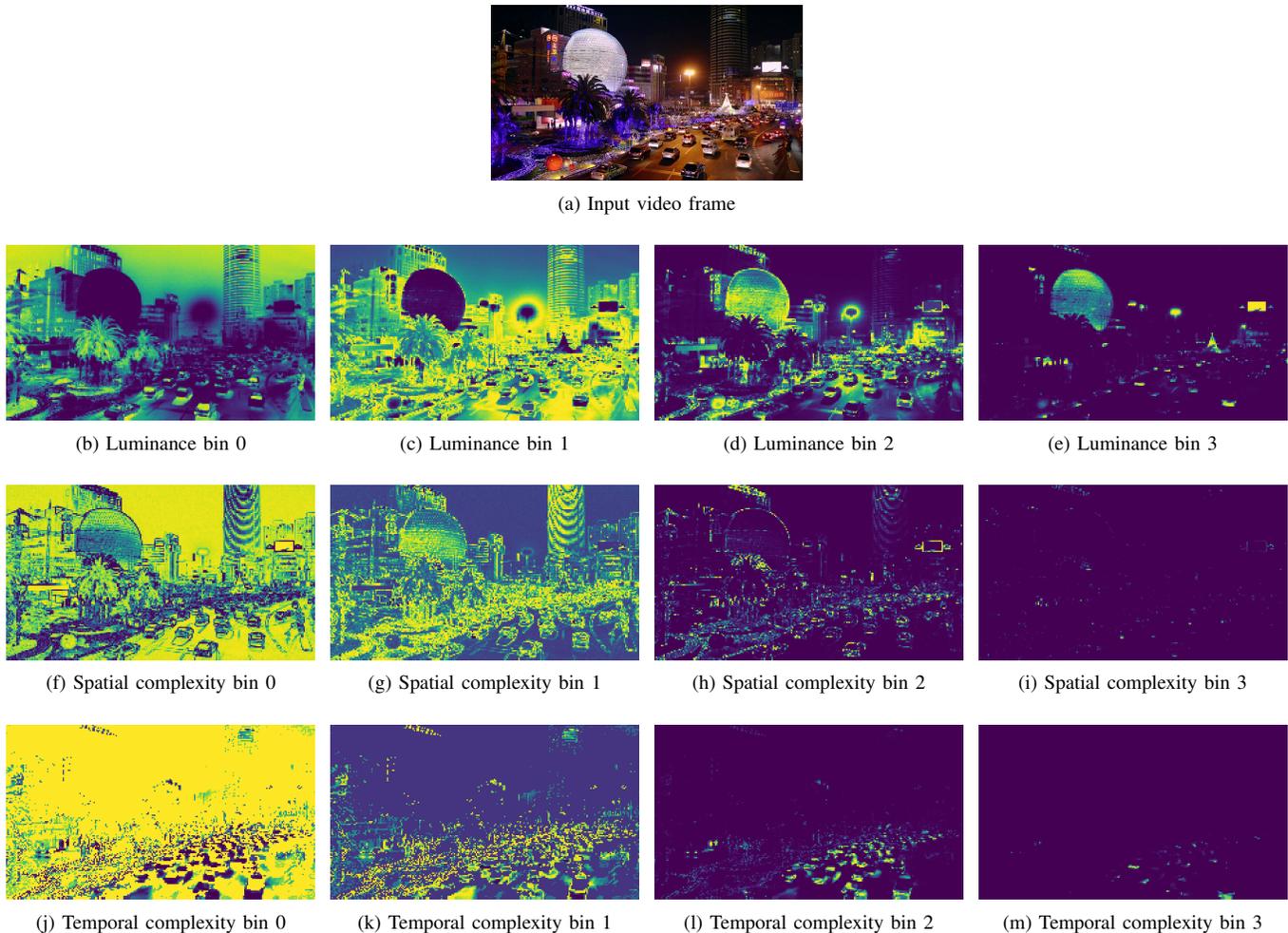


Fig. 3. A sample video frame and the soft-classification of 8×8 cuts into four bins based on luminance, spatial, and temporal complexity features.

pooling to be numerically unstable for low quality inputs. Hence, SSIM is aggregated as

$$SSIM_b^{(L)} = 1 - \left(\frac{\sum_c M_b^{(L)}(c) \times (1 - SSIM(c))^3}{\sum_c M_b^{(L)}(c)} \right)^{1/3}. \quad (54)$$

Prior work on spatial quality aggregation [43] has revealed that low-quality spatial regions are more salient to observers and, therefore, largely determine the overall perceived quality of the input stimulus. In line with this observation, we posit that the most salient bin of each type, which corresponds to spatial regions in the image, is that which has the lowest quality. So, we aggregate bin-level quality features into a single feature per scale by using the quality of the worst-quality bin. For quality features such as SSIM and VIF (i.e., a greater value implies better quality), this involves applying a min operation over bin-level features, while for distortion features such as SRRED and FOSD (i.e., a greater value implies worse quality), this involves applying a max operation over bin-level features.

Although perceptually motivated, using only the quality features from the worst image regions ignores other regions of the input frame. So, we also obtain global estimates of quality-aware features at each scale by computing “unweighted

averages,” i.e., the mean value without binning or weighting. These features also function as a baseline against which the efficacy of binned-weighting is evaluated. The results of detailed ablation studies studying the effect of weighted aggregation and the various bin types are presented in Section IV-C.

The procedure described thus yields binned-weighted and unweighted quality features computed at four scales. Finally, these features are aggregated across the four scales to yield a single multi-scale feature using the multi-scale fusion weights proposed in MS-SSIM [47]. Let w_s denote the MS-SSIM fusion weight for scale s and let a quality feature at that scale be Q_s . Then, the multi-scale quality feature was computed as:

$$Q = \frac{\sum_s w_s Q_s}{\sum_s w_s}. \quad (55)$$

The combination of global NSS features described in Section III-E and the multi-scale features computed as described above yields a total of 232 features.

IV. EVALUATING OBJECTIVE QUALITY MODELS

In this section, we report the results of evaluating the effectiveness of the Cut-FUNQUE Model against a set of baseline quality models from the literature. In addition, we

analyzed various combinations of the novel proposals made in the Cut-FUNQUE model in the form of an ablation study. To conduct these evaluations, we utilized the recently developed LIVE-TMHDR database [48].

A. The LIVE-TMHDR Database

The LIVE-TMHDR Database is a recently-developed first-of-its-kind large-scale subjective database containing a set of 40 source HDR videos, 20 of which were created by professional studios and 20 generated by amateur iPhone users. The set of 40 source contents were subjected to both tone-mapping, using a wide variety of algorithms and manually tone-mapped by an expert colorist, and video compression, using the x264 [49] encoder, yielding a total of 15,000 test contents.

The new database represents twelve tone-mapping operators (TMOs), of which ten were made open-source. The ten open-source TMOs and their key features may be summarized as follows.

- 1) **Hable** [50] - A parameter-free pointwise non-linear transform originally designed for use in the video game *Uncharted 2*.
- 2) **Reinhard02** [51] - A point non-linearity to map luminances from HDR to SDR.
- 3) **Durand02** [52] - Uses a “fast bilateral filter” to decompose the luminances of HDR frames into “base” and “detail” layers.
- 4) **Shan12** [53] - Uses an edge-aware stationary wavelet transform (SWT) [54].
- 5) **Reinhard12** [55] - Uses color-appearance models applied in a local manner.
- 6) **Eilertsen15** [56] - Applies a “fast detail extraction” method to obtain a base-detail decomposition and applies a dynamic tone-curve.
- 7) **Oskarsson17** [57] - Uses Dynamic Programming to cluster values in the input image channels.
- 8) **Rana19** [58] - Uses a Generative Adversarial Network (GAN) to create a fully-convolutional, parameter-free TMO.
- 9) **Yang21** [59] - Uses a deep convolutional neural network (CNN) to transform a multi-scale Laplacian pyramid decomposition of each input HDR frame.
- 10) **ITU21** [60] - A parameter-free TMO proposed by the ITU in Recommendation BT.2446 (“Approach A”).

In addition to these methods, the DolbyVision TMO (DV), created by Dolby as part of the DolbyVision HDR standard, and the Color Space Transform (CST) provided in DaVinci Resolve, which is a popular gamut/tone-mapping tool used by colorists, were also included in the database. Finally, all videos were encoded at three compression levels using the libx264 [49] encoder. Due to its large size and diversity of source contents and distortions, we chose LIVE-TMHDR as a test bench for evaluating quality prediction models aimed at the practical delivery of tone-mapped HDR videos¹.

¹Experiments using the LIVE-TMHDR database were conducted at the University of Texas at Austin by university-affiliated authors.

B. Evaluation Protocol

To evaluate Cut-FUNQUE and compare it with quality prediction models from the literature, we conducted content-separated cross-validation and report the median Pearson Correlation Coefficient (PCC), Spearman’s Rank Order Correlation Coefficient (SROCC), and Root Mean Square Error (RMSE) over 100 random splits of the LIVE-TMHDR database.

In each split, 80% of the data was used for training and 20% was used for testing. When generating random splits, we ensured that test contents from the same source HDR video were not present in both the training and test splits. To show each quality prediction feature set in its best light, we experimented with three regression models (wherever feasible) - the Linear Support Vector Regressor (SVR), Gaussian SVR, and Random Forest Regressor - to map features to quality scores. Hyper-parameters of each regressor were also tuned using cross-validation, and the regressor that achieved the highest cross-validation accuracy was selected. In the case of MSML, the partial least-squares projection matrix was recalibrated on the training dataset of each cross-validation split for a fair evaluation.

C. Results of Quality Prediction on LIVE-TMHDR

To illustrate the efficacy of Cut-FUNQUE, we evaluated its performance against 15 TM-HDR quality prediction models from the literature. TMQI [16], FSITM [21], and FFTMI [23] are full-reference image quality assessment (FR IQA) models, TMVQI [25] and FUNQUE+ [14] are FR video quality assessment (FR VQA) models, and BRISQUE [19], NIQE [20], DIIVINE [61], BTMQI [26], RcNet [28], HIGRADE [27], and MSML [29] are no-reference (NR) IQA models. We adapted all IQA models to videos by applying them framewise.

The median cross-validation accuracy achieved by each quality prediction model on the LIVE-TMHDR database is presented in Table I. From Table I, it may be seen that Cut-FUNQUE significantly outperformed, by over 15%, nearly all of the compared quality models on the LIVE-TMHDR database in terms of quality prediction accuracy. The only existing quality model rivaling Cut-FUNQUE is the deep MSML model that computes a set of 9216 features from a pre-trained ResNet-50 model. By contrast, Cut-FUNQUE relies on an efficient Haar wavelet transform and computation sharing between features.

To formally evaluate the efficiency of Cut-FUNQUE, we used the Performance Application Programming Interface (PAPI) [62] to count the number of Floating Point Operations (FLOPs) required by each quality model evaluated in Table I. This measurement provides a theoretical estimate of each model’s computational complexity. In addition, we also measured the running time of each quality model over 50 frames of a 1080p input video pair on a six-core Intel Core i7-8700 CPU, which has a clock frequency of 3.20GHz. Using this measurement, we report the average time taken per frame by each quality model as a descriptor of its practical computational complexity. These measurements, in units of Giga FLOPs and seconds, are also included in Table I. From

TABLE I
EVALUATION OF QUALITY PREDICTION MODELS IN TERMS OF MEDIAN CROSS-VALIDATION ACCURACY AND COMPUTATIONAL COMPLEXITY

Model	Regressor	PCC	SROCC	RMSE	GFLOPs/Frame	Time (Sec/Frame)
Y-FUNQUE+ [14]	RandomForest	0.4524	0.4343	9.4352	0.1042	2.414
BTMQI [26]	GaussianSVR	0.4705	0.4663	9.2238	0.1199	4.612
FSITM [21]	LinearSVR	0.4813	0.4626	8.9212	8.9487	4.930
NIQE [20]	GaussianSVR	0.4805	0.4746	9.5563	2.3654	3.198
BRISQUE [19]	LinearSVR	0.4811	0.4833	8.9869	0.2120	0.893
DIIVINE [61]	GaussianSVR	0.4794	0.4925	9.2879	20.8771	155.373
TMQI [16]	GaussianSVR	0.5062	0.4956	8.6897	0.9061	2.374
FUNQUE [13]	RandomForest	0.5082	0.4949	8.8863	0.1716	2.471
TMVQI [25]	RandomForest	0.5198	0.4969	8.8697	1.4164	3.002
FFTMI [23]	GaussianSVR	0.5298	0.5315	8.8559	27.5161	14.407
3C-FUNQUE+ [14]	RandomForest	0.5817	0.5661	8.6568	0.3667	2.825
ReNet [28]	Random Forest	0.5985	0.5824	8.2417	134.5597	1280.730
HIGRADE [27]	GaussianSVR	0.6682	0.6698	8.2619	2.6533	3.709
MSML [29]	Linear SVR	0.7883	0.7740	6.8090	67.2578	412.438
Cut-FUNQUE	Random Forest	0.7783	0.7781	6.4187	2.9257	6.343

TABLE II
QUALITY PREDICTION ACCURACY OF CUT-FUNQUE VARIANTS USING VARIOUS ENCODING FUNCTIONS

Encoding Function	PCC	SROCC	RMSE
PUColor	0.7783	0.7781	6.4187
PU21	0.7765	0.7686	6.4313
PQ	0.7677	0.7619	6.4738

TABLE III
EFFECT OF PROGRESSIVELY REMOVING WEIGHTED-AGGREGATED FEATURES FROM CUT-FUNQUE

Model	PCC	SROCC	RMSE
Cut-FUNQUE	0.7783	0.7781	6.4187
w/o T-Weighted Features	0.7636	0.7540	6.4808
w/o L-weighted Features	0.7513	0.7520	6.6442
w/o S-weighted Features	0.7389	0.7337	6.9726

this, it may be seen that Cut-FUNQUE achieves comparable accuracy as MSML, while being $\sim 23\times$ as efficient, in terms of FLOPs.

To analyze the novel proposals in Cut-FUNQUE, we conducted an ablation study by comparing variants of Cut-FUNQUE that use different encoding functions (in Table II) and subsets of low-level cut weighted features (luminance-, spatial-, and temporal-complexity) (in Table III). From these results, it may be seen that replacing PUColor with PQ or PU21, or removing weighted-aggregated features decreases accuracy. These results validate the design choices made in Cut-FUNQUE.

Finally, we lend further weight to the comparison of quality prediction accuracy on LIVE-TMHDR by conducting a statistical significance test on the differences in accuracy. Specifically, we conducted one-sided Welch’s t-tests to evaluate the statistical significance of the observed differences in prediction accuracies. A one-sided Welch’s t-test was preferred over a traditional Student’s t-test, since a Welch’s t-test accounts for unequal population variances [63]. Table IV presents the results of pairwise statistical significance comparisons. An entry of “1” (“0”) denotes that the quality model in the row achieved statistically significantly superior (inferior) accuracy compared to the quality model in the column. An entry of “-” denotes that the differences are not statistically significant. From this table,

it may be seen that MSML and Cut-FUNQUE outperformed all other quality models, and that the difference between the two top performers was not statistically significant.

V. CONCLUSION

We have developed and validated a novel model called Cut-FUNQUE, which has been designed to predict the perceptual quality of tone-mapped and compressed HDR videos. An analysis of the quality prediction accuracy and computational complexity of Cut-FUNQUE demonstrates that it achieves comparable accuracy as the SOTA deep quality model MSML, while using a fraction (less than 5%) of its computational cost.

Despite its benefits, Cut-FUNQUE could be improved further by applying more accurate global and local color and contrast quality features. This work also leaves room for studies of the trade-offs between reproducing the appearance of a source HDR video and the standalone quality of corresponding tone-mapped videos. Finally, the development of domain-specific deep neural networks that leverage knowledge of tone-mapping distortions may lead to further improvements in quality prediction accuracy. The primary challenge facing such approaches is the scarcity of publicly available HDR videos and even larger tone mapped video quality databases.

REFERENCES

- [1] IEC, “Multimedia systems and equipment - colour measurement and management - part 2-1: Colour management - default RGB colour space - sRGB,” 1999.
- [2] ITU-R, “ITU-R BT.709: Parameter values for the hdtv standards for production and international programme exchange,” 2011.
- [3] —, “ITU-R BT.2100: Image parameter values for high dynamic range television for use in production and international programme exchange,” 2018.
- [4] S. Standard, “High dynamic range electro-optical transfer function of mastering reference displays,” *SMPTE ST*, vol. 2084, no. 2014, p. 11, 2014.
- [5] T. Borer and A. Cotton, “A display-independent high dynamic range television system,” *SMPTE Motion Imaging Journal*, vol. 125, no. 4, pp. 50–56, 2016.
- [6] CTA. Television technology consumer definitions. [Online]. Available: <https://cdn.cta.tech/cta/media/media/membership/pdfs/videotechnology-consumer-definitions.pdf>
- [7] HDR10+ Technologies, LLC. (2019) HDR10+ system whitepaper. [Online]. Available: https://hdr10plus.org/wp-content/uploads/2019/08/HDR10_WhitePaper.pdf

TABLE IV

RESULTS OF WELCH'S T-TEST APPLIED TO PAIRS OF OBJECTIVE QUALITY MODELS.

AN ENTRY OF "1" ("0") DENOTES THAT THE MODEL IN THE ROW ACHIEVES SUPERIOR (INFERIOR) ACCURACY TO THE MODEL IN THE COLUMN, WHILE AN ENTRY OF "-" DENOTES THAT THE DIFFERENCE IN ACCURACY IS NOT STATISTICALLY SIGNIFICANT.

	Y-FUNQUE+	BTMQI	FSITM	NIQE	BRISQUE	DIIVINE	TMQI	FUNQUE	TMVQI	FFTMI	3C-FUNQUE+	ReNet	HIGRADE	MSML	Cut-FUNQUE
Y-FUNQUE+	-	-	-	-	0	0	0	0	0	0	0	0	0	0	0
BTMQI	-	-	-	-	0	0	0	0	0	0	0	0	0	0	0
FSITM	-	-	-	-	-	-	0	0	0	0	0	0	0	0	0
NIQE	-	-	-	-	-	0	0	0	0	0	0	0	0	0	0
BRISQUE	1	1	-	-	-	-	-	-	-	-	-	0	0	0	0
DIIVINE	1	1	-	-	-	-	-	-	0	-	0	0	0	0	0
TMQI	1	1	1	1	-	-	-	-	-	-	0	0	0	0	0
FUNQUE	1	1	1	1	-	-	-	-	0	-	0	0	0	0	0
TMVQI	1	1	1	1	-	-	-	-	0	-	0	0	0	0	0
FFTMI	1	1	1	1	1	1	-	1	1	-	-	0	0	0	0
3C-FUNQUE+	1	1	1	1	-	-	-	-	-	-	-	0	0	0	0
ReNet	1	1	1	1	1	1	1	1	1	1	-	0	0	0	0
HIGRADE	1	1	1	1	1	1	1	1	1	1	1	1	-	0	0
MSML	1	1	1	1	1	1	1	1	1	1	1	1	1	-	-
Cut-FUNQUE	1	1	1	1	1	1	1	1	1	1	1	1	1	-	-

[8] C. Forrester, "SkyPerfect offers UHD-HDR by DTH." [Online]. Available: <https://advanced-television.com/2015/11/04/skyperfect-offers-uhd-hdr-by-dth/>

[9] (2016) An introduction to Dolby Vision. [Online]. Available: https://professional.dolby.com/siteassets/pdfs/dolby-vision-whitepaper_an-introduction-to-dolby-vision_0916.pdf

[10] CNET. Best TV for 2024: We Tested Samsung, LG, TCL, Vizio and More. [Online]. Available: <https://www.cnet.com/tech/home-entertainment/best-tv/>

[11] M. Čadík, M. Wimmer, L. Neumann, and A. Artusi, "Evaluation of HDR tone mapping methods using essential perceptual attributes," *Computers & Graphics*, vol. 32, no. 3, pp. 330–349, 2008.

[12] J. Morovic and M. R. Luo, "The fundamentals of gamut mapping: A survey," *Journal of Imaging Science and Technology*, vol. 45, no. 3, pp. 283–290, 2001.

[13] A. K. Venkataramanan, C. Stejerean, and A. C. Bovik, "FUNQUE: Fusion of unified quality evaluators," in *2022 IEEE International Conference on Image Processing (ICIP)*, 2022, pp. 2147–2151.

[14] A. K. Venkataramanan, C. Stejerean, I. Katsavounidis, and A. C. Bovik, "One transform to compute them all: Efficient fusion-based full-reference video quality assessment," *IEEE Transactions on Image Processing*, vol. 33, pp. 509–524, 2024.

[15] T. O. Aydin, R. Mantiuk, K. Myszkowski, and H.-P. Seidel, "Dynamic range independent image quality assessment," *ACM Trans. Graph.*, vol. 27, no. 3, p. 1–10, aug 2008. [Online]. Available: <https://doi.org/10.1145/1360612.1360668>

[16] H. Yeganeh and Z. Wang, "Objective quality assessment of tone-mapped images," *IEEE Transactions on Image Processing*, vol. 22, no. 2, pp. 657–667, 2013.

[17] Z. Wang, A. Bovik, H. Sheikh, and E. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE Transactions on Image Processing*, vol. 13, no. 4, pp. 600–612, 2004.

[18] D. Kundu and B. L. Evans, "Visual attention guided quality assessment of tone-mapped images using scene statistics," in *2016 IEEE International Conference on Image Processing (ICIP)*, 2016, pp. 96–100.

[19] A. Mittal, A. K. Moorthy, and A. C. Bovik, "No-reference image quality assessment in the spatial domain," *IEEE Transactions on Image Processing*, vol. 21, no. 12, pp. 4695–4708, 2012.

[20] A. Mittal, R. Soundararajan, and A. C. Bovik, "Making a "completely blind" image quality analyzer," *IEEE Signal Processing Letters*, vol. 20, no. 3, pp. 209–212, 2013.

[21] H. Ziaei Nafchi, A. Shahkolaei, R. Farrahi Moghaddam, and M. Cheriet, "FSITM: A feature similarity index for tone-mapped images," *IEEE Signal Processing Letters*, vol. 22, no. 8, pp. 1026–1029, 2015.

[22] L. Zhang, L. Zhang, X. Mou, and D. Zhang, "FSIM: A feature similarity index for image quality assessment," *IEEE Transactions on Image Processing*, vol. 20, no. 8, pp. 2378–2386, 2011.

[23] L. Krasula, K. Fliegel, and P. Le Callet, "FFTMI: Features fusion for natural tone-mapped images quality evaluation," *IEEE Transactions on Multimedia*, vol. 22, no. 8, pp. 2038–2047, 2020.

[24] L. Krasula, M. Narwaria, K. Fliegel, and P. L. Callet, "Rendering of HDR content on LDR displays: an objective approach," in *Applications of Digital Image Processing XXXVIII*, A. G. Tescher, Ed., vol. 9599, International Society for Optics and Photonics. SPIE, 2015, p. 95990X. [Online]. Available: <https://doi.org/10.1117/12.2186388>

[25] H. Yeganeh, S. Wang, K. Zeng, M. Eisapour, and Z. Wang, "Objective quality assessment of tone-mapped videos," in *2016 IEEE International Conference on Image Processing (ICIP)*, 2016, pp. 899–903.

[26] K. Gu, S. Wang, G. Zhai, S. Ma, X. Yang, W. Lin, W. Zhang, and W. Gao, "Blind quality assessment of tone-mapped images via analysis of information, naturalness, and structure," *IEEE Transactions on Multimedia*, vol. 18, no. 3, pp. 432–443, 2016.

[27] D. Kundu, D. Ghadiyaram, A. C. Bovik, and B. Evans, "No-reference quality assessment of tone-mapped HDR pictures," *IEEE Transactions on Image Processing*, vol. 26, no. 6, pp. 2957–2971, 2017.

[28] C. S. Ravuri, R. Sureddi, S. V. R. Dendi, S. Raman, and S. S. Channappayya, "Deep no-reference tone mapped image quality assessment," in *2019 53rd Asilomar Conference on Signals, Systems, and Computers*. IEEE, 2019, pp. 1906–1910.

[29] Q. He, D. Li, T. Jiang, and M. Jiang, "Quality assessment for tone-mapped HDR images using multi-scale and multi-layer information," in *2018 IEEE International Conference on Multimedia & Expo Workshops (ICMEW)*, 2018, pp. 1–6.

[30] M. A. Saad, A. C. Bovik, and C. Charrier, "Blind prediction of natural video quality," *IEEE Transactions on Image Processing*, vol. 23, no. 3, pp. 1352–1365, 2014.

[31] Z. Tu, Y. Wang, N. Birkbeck, B. Adsumilli, and A. C. Bovik, "UGC-VQA: Benchmarking blind video quality assessment for user generated content," *IEEE Transactions on Image Processing*, vol. 30, pp. 4449–4464, 2021.

[32] Z. Li, A. Aaron, I. Katsavounidis, A. Moorthy, and M. Manohara, "Toward a practical perceptual video quality metric," *The Netflix Tech Blog*, vol. 6, p. 2, 2016.

[33] A. K. Venkataramanan, C. Stejerean, I. Katsavounidis, and A. C. Bovik, "A funque approach to the quality assessment of compressed hdr videos," *arXiv preprint arXiv:2312.08524*, 2023.

[34] R. K. Mantiuk and M. Azimi, "PU21: A novel perceptually uniform encoding for adapting existing quality metrics for HDR," in *2021 Picture Coding Symposium (PCS)*, 2021, pp. 1–5.

[35] J. P. Ebenezer, Z. Shang, Y. Wu, H. Wei, S. Sethuraman, and A. C. Bovik, "Making video quality assessment models robust to bit depth," *IEEE Signal Processing Letters*, vol. 30, pp. 488–492, 2023.

[36] S. Li, F. Zhang, L. Ma, and K. N. Ngan, "Image quality assessment by separately evaluating detail losses and additive impairments," *IEEE Transactions on Multimedia*, vol. 13, no. 5, pp. 935–949, 2011.

[37] R. Soundararajan and A. C. Bovik, "Video quality assessment by reduced reference spatio-temporal entropic differencing," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 23, no. 4, pp. 684–694, 2013.

[38] R. K. Mantiuk, M. Kim, M. Ashraf, Q. Xu, M. R. Luo, J. Martinovic, and S. Wuerger, "Practical color contrast sensitivity functions for luminance levels up to 10000 cd/m²," *Color and Imaging Conference*, vol. 2020, no. 28, pp. 1–6, 2020. [Online]. Available: <https://www.ingentaconnect.com/content/ist/cic/2020/00002020/00000028/art00002>

[39] A. M. Derrington, J. Krauskopf, and P. Lennie, "Chromatic mechanisms in lateral geniculate nucleus of macaque." *The Journal of physiology*, vol. 357, no. 1, pp. 241–265, 1984.

[40] J. Schanda, *Colorimetry: understanding the CIE system*. John Wiley & Sons, 2007.

[41] K. Gu, M. Liu, G. Zhai, X. Yang, and W. Zhang, "Quality assessment considering viewing distance and image resolution," *IEEE Transactions on Broadcasting*, vol. 61, no. 3, pp. 520–531, 2015.

[42] A. B. Watson, G. Y. Yang, J. A. Solomon, and J. Villasenor, "Visibility of wavelet quantization noise," *IEEE Transactions on Image Processing*, vol. 6, no. 8, pp. 1164–1175, 1997.

- [43] A. K. Venkataramanan, C. Wu, A. C. Bovik, I. Katsavounidis, and Z. Shahid, "A hitchhiker's guide to structural similarity," *IEEE Access*, vol. 9, pp. 28 872–28 896, 2021.
- [44] H. R. Sheikh and A. C. Bovik, "Image information and visual quality," *IEEE Transactions on image processing*, vol. 15, no. 2, pp. 430–444, 2006.
- [45] K. Sharifi and A. Leon-Garcia, "Estimation of shape parameter for generalized gaussian distributions in subband decompositions of video," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 5, no. 1, pp. 52–56, 1995.
- [46] N.-E. Lasmaz, Y. Stitou, and Y. Berthoumieu, "Multiscale skewed heavy tailed model for texture analysis," in *2009 16th IEEE International Conference on Image Processing (ICIP)*, 2009, pp. 2281–2284.
- [47] Z. Wang, E. Simoncelli, and A. Bovik, "Multiscale structural similarity for image quality assessment," in *The Thirty-Seventh Asilomar Conference on Signals, Systems & Computers, 2003*, vol. 2, 2003, pp. 1398–1402 Vol.2.
- [48] A. K. Venkataramanan and A. C. Bovik, "Subjective quality assessment of compressed tone-mapped high dynamic range videos," *Manuscript Under Preparation*, vol. 1, 2024.
- [49] VideoLAN, "x264." [Online]. Available: <https://code.videolan.org/videolan/x264.git>
- [50] J. Hable. Uncharted 2: HDR lighting. [Online]. Available: <https://www.gdcvault.com/play/1012351/Uncharted-2-HDR>
- [51] E. Reinhard, M. Stark, P. Shirley, and J. Ferwerda, "Photographic tone reproduction for digital images," in *Proceedings of the 29th Annual Conference on Computer Graphics and Interactive Techniques*, ser. SIGGRAPH '02. Association for Computing Machinery, 2002, p. 267–276.
- [52] F. Durand and J. Dorsey, "Fast bilateral filtering for the display of high-dynamic-range images," in *Proceedings of the 29th Annual Conference on Computer Graphics and Interactive Techniques*, ser. SIGGRAPH '02. New York, NY, USA: Association for Computing Machinery, 2002, p. 257–266.
- [53] Q. Shan, T. DeRose, and J. Anderson, "Tone mapping high dynamic range videos using wavelets," *Pixar Technical Memo*, 2012.
- [54] G. P. Nason and B. W. Silverman, *The Stationary Wavelet Transform and some Statistical Applications*. Springer New York, 1995, pp. 281–299.
- [55] E. Reinhard, T. Pouli, T. Kunkel, B. Long, A. Ballestad, and G. Damborg, "Calibrated Image Appearance Reproduction," *ACM Trans. Graph.*, vol. 31, no. 6, Nov 2012.
- [56] G. Eilertsen, R. K. Mantiuk, and J. Unger, "Real-time noise-aware tone mapping," *ACM Trans. Graph.*, vol. 34, no. 6, Nov 2015.
- [57] M. Oskarsson, "Temporally consistent tone mapping of images and video using optimal k-means clustering," *Journal of Mathematical Imaging and Vision*, vol. 57, no. 2, pp. 225–238, Feb 2017.
- [58] A. Rana, P. Singh, G. Valenzise, F. Dufaux, N. Komodakis, and A. Smolic, "Deep tone mapping operator for high dynamic range images," *IEEE Transactions on Image Processing*, vol. 29, pp. 1285–1298, 2020.
- [59] J. Yang, Z. Liu, M. Lin, S. Yanushkevich, and O. Yadid-Pecht, "Deep reformulated laplacian tone mapping," *arXiv preprint arXiv:2102.00348*, 2021.
- [60] ITU-R, "ITU-R BT.2446: Methods for conversion of high dynamic range content to standard dynamic range content and vice-versa," 2021.
- [61] A. K. Moorthy and A. C. Bovik, "Blind image quality assessment: From natural scene statistics to perceptual quality," *IEEE Transactions on Image Processing*, vol. 20, no. 12, pp. 3350–3364, 2011.
- [62] S. Browne, J. Dongarra, N. Garner, G. Ho, and P. Mucci, "A portable programming interface for performance evaluation on modern processors," *The International Journal of High Performance Computing Applications*, vol. 14, no. 3, pp. 189–204, 2000. [Online]. Available: <https://doi.org/10.1177/109434200001400303>
- [63] M. Delacre, D. Lakens, and C. Leys, "Why psychologists should by default use Welch's t-test instead of Student's t-test," *International Review of Social Psychology*, vol. 30, no. 1, pp. 92–101, 2017.