**Contact**    maabed@gmail.com OR {gukyeong.kwon,alregib}@gatech.edu
`http://ghassanalregib.com/`

# POWER OF TEMPOSPATIALLY UNIFIED SPECTRAL DENSITY FOR PERCEPTUAL VIDEO QUALITY ASSESSMENT

*Mohammed A. Aabed, Gukyeong Kwon, and Ghassan AlRegib*

Center for Signal and Information Processing (CSIP)
School of Electrical and Computer Engineering
Georgia Institute of Technology Atlanta, Georgia 30332, U.S.A.
{maabed,gukyeong.kwon,alregib}@gatech.edu

## ABSTRACT

We propose a perceptual video quality assessment (PVQA) metric for distorted videos by analyzing the power spectral density (PSD) of a group of pictures. This is an estimation approach that relies on the changes in video dynamic calculated in the frequency domain and are primarily caused by distortion. We obtain a feature map by processing a 3D PSD tensor obtained from a set of distorted frames. This is a full-reference tempospatial approach that considers both temporal and spatial PSD characteristics. This makes it ubiquitously suitable for videos with varying motion patterns and spatial contents. Our technique does not make any assumptions on the coding conditions, streaming conditions or distortion. This approach is also computationally inexpensive which makes it feasible for real-time and practical implementations. We validate our proposed metric by testing it on a variety of distorted sequences from PVQA databases. The results show that our metric estimates the perceptual quality at the sequence level accurately. We report the correlation coefficients with the differential mean opinion scores (DMOS) reported in the databases. The results show high and competitive correlations compared with the state of the art techniques.

*Index Terms*— Perceptual quality, video quality, perception, human visual system, 3D power spectral density

## 1. INTRODUCTION

The proliferation of visual media, in general, and video streaming services and applications, in particular, in recent years has increased the need for efficient communication, bandwidth and streaming. Video traffic in 2015 accounted for over 55% of global IP traffic. By 2020, it is estimated that a growth of 68% in global mobile connections will occur reaching 11.6 billion mobile connections. Mobile video traffic will account for over 75% of that total. In fact, it will take an individual five million years to watch the amount of monthly video traffic transmitted through global IP networks [1]. Furthermore, the development and enhancements of video coding standards have been very active over the past decade. In addition to the release of H.265/MPEG-H Part 2 High Efficiency Video Coding (HEVC) in 2013, several development activities from industrial corporations emerged outside the umbrella of Moving Picture Experts Group (MPEG) and International Telecommunications Union (ITU). The recently formed Alliance for Open Media (AOMedia) includes several major industry leaders whose main purpose is developing a true universal royalty-free video coding standard. The alliance is anticipating the release of its first standard in 2017, AOMedia Video 1 (AV1) [2, 3]. This also coincides with the Chinese government and companies' ongoing efforts towards AVS2 [4]. Nonetheless, video coding development focuses on developing a standard format for the bitstream and decoder mainly. The process involves describing general coding tools without explicitly defining their design. This flexible standardization procedure leaves room for optimizations and innovation but *comes with no guarantees of perceptual video quality*. Henceforth, the importance of quality of experience (QoE) has been critically emphasized in this domain.
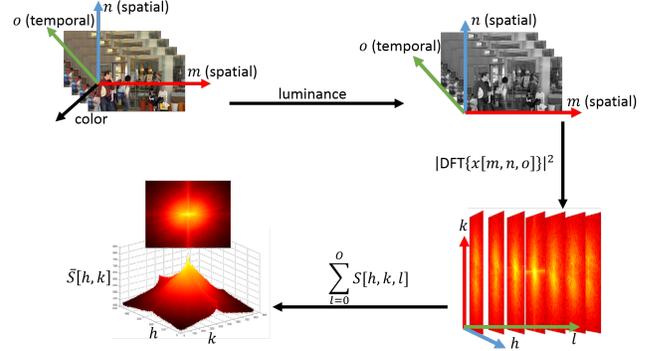
To establish stable video operations and services while maintaining high quality of experience, *perceptual* video quality assessment (PVQA) becomes an essential research topic in video technology. A survey published in 2015 revealed that one out of five viewers will abandon a poor streaming service immediately. Furthermore, 75% of the users will tolerate a bad stream for up to four minutes before switching to a more reliable one [5]. The significance of PVQA is not limited to quality control only. Perceptual video quality plays a pivotal role in designing and improving super resolution and video enhancements algorithms. PVQA is also conjointly associated with the evolving understanding of the human visual system (HVS) and visual perception in the computational neuroscience community. The two research domains complement one another filling the gaps in our understanding, processing and development of visual media technology. Hence, this paper addresses this problem and introduces a new framework for video quality assessment.

Several video quality assessment approaches have been proposed over the past decade. PVQA has several chal-

lenges including the incorporating visual perception characteristics, feature selection and crafting, distortion detection and tracking, and pooling optimization among others. Tempospatial feature processing has been investigated and proposed in several ways in past works. In [6], the authors use tempospatial Gabor filters and motion trajectory to evaluate spatial, temporal and tempospatial quality of the videos. The work in [7] proposed a hierarchal statistical processing model for video quality monitoring using pixel-level optical flow motion fields. Soundararajan and Bovik [8] utilize natural video statistics in the wavelet domain and entropic differences to predict video quality. Saad *et. al* [9] proposes a no-reference video quality measure relying on tempospatial natural statistics and motion models using discrete cosine transform (DCT). 3D shearlet transform is applied to videos to capture directions of curvlinear singularities and anisotropic features in [10]. The authors in [11] introduce 3D singular value decomposition as content based transform and measure the quality of video by comparing singular values of original and distorted videos. In [12], a no-reference video quality models of intrinsic statistical regularities observed in natural videos, which are used to quantify distortions. This work introduces a new perceptual video quality framework using tempospatial power spectral density (PSD) processing. To the best of the authors knowledge, our work is the first to explore and utilize PSD in video quality assessment.

In this paper, we propose a new approach to predict the quality of video through the analysis on 3D PSD. We propose a full-reference perceptual video quality metric based on the disruptions in the power of tempospatially unified spectra. Our approach characterizes distortions through the statistical features in 3D PSD by fusing tempospatial power spectral density (TPSD) planes of the distorted and pristine videos to estimate the perceived video quality. The PSD is one of the distinctive frequency domain characteristics of a signal. In addition to the power distribution of video frames, PSD also thoroughly captures scene features and objects [13]. Through the 3D processing in our metric, spatial and temporal distortions are analyzed simultaneously. The combined effect from both distortions is effectively captured in the same framework. Moreover, the discrete Fourier transform (DFT) is the main operation required to calculate the PSD, which is a very computationally simple domain transform compared to other operations such as wavelet and curvelet transforms. The computational simplicity enables real-time processing, future extensions and application diversity.

The rest of this paper is organized as follow. We explain statistical features in 3D TPSD and processing flow of our propose method in Section 2. In Section 3, we validate our metric by examining the correlations with the human mean opinion scores (MOS). We also compare our proposed metric against well-known and state of the art VQA metrics. Section 4 concludes the paper and highlights future directions.



**Fig. 1**. 3D power spectral density tensor-level processing flowchart.

## 2. PROPOSED METHOD

### 2.1. 3D Power Spectral Density

A 3D discrete time-space video signal is defined as $x[m, n, o] \in \mathscr{R}^{M \times N \times O}$, with the one grayscale (luma) channel, where $m$ and $n$ are the spacial indices of the 2D frame and $o$ is the temporal (frame) index. The frequency response of 3D discrete time signal $x[m, n, o]$ is derived by calculating 3D DFT, $X[h, k, l] \in \mathscr{C}^{M \times N \times O}$. The 3D discrete PSD, $\mathcal{S}[h, k, l] \in \mathscr{R}^{M \times N \times O}$, can be estimated using Parseval's theorem as follows:

$$\mathcal{S}[h, k, l] = \frac{1}{MNO} |X[h, k, l]|^2, \qquad (1)$$

where $h$, $k$ and $l$ are the discrete frequency indices.

In order to calculate the 2D time-aggregated tempospatial PSD (TPSD) plane at every spatial frequency, $\overline{\mathcal{S}}[h, k]$, the expression in (1) is integrated over the temporal axis, $O$. That is,
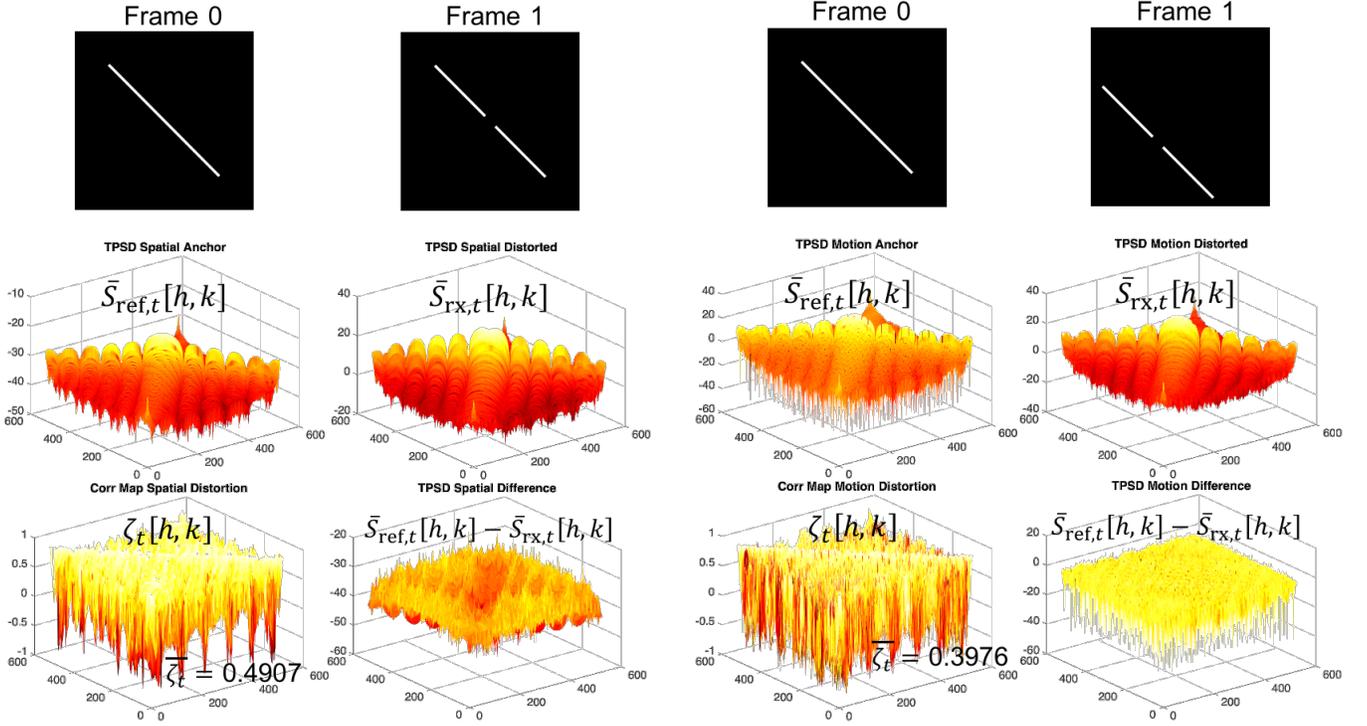
$$\overline{\mathcal{S}}[h, k] = \sum_{l=0}^{O} \mathcal{S}[h, k, l], \qquad (2)$$

where $\overline{\mathcal{S}}[h, k] \in \mathscr{R}^{M \times N}$. Figure 1 illustrates the processing framework for a tensor of frames of size $M \times N \times O$.

### 2.2. Video Quality Estimation based on 3D PSD

Distortions in a video change the tempospatial characteristics in the pixel domain causing a distribution of the signal's PSD. These changes can be captured using 2D time-aggregated TPSD plane, $\overline{\mathcal{S}}[h, k]$. The deviation of the distorted TPSD from the original free of distortion can be captured in several ways to reflect the change in the energy field. We estimate this variability and map it to perception by measuring a local cross-correlations map between the distorted and anchor TPSD.

Two videos, an anchor video free of distortion and a distorted video, are divided into a set of tensors. For simplicity,

(a) Temporal distortion without motion. Frame 1 in the anchor video is identical to Frame 0.

(b) Temporal distortion with motion. The line in anchor Frame 1 is shifted downwards from its original location in Frame 0 to intorduce a tempospatail change.

**Fig. 2**. Simple examples illustrating the principles underlying the proposed metric. Both examples are composed of two frames only ($O = 2$) where we show the distorted versions of the Frame 1. For every sequence, we show the the anchor and distorted tempospatial planes, difference map and local-cross correlation distortion map.

we assume the tensors to be of equal size $M \times N \times O$. In practice, tensors sizes may vary depending on coding group of pictures, scene boundaries, processing efficiency, etc. Let the $t^{\text{th}}$ tensor in the anchor and distorted video be denoted as $x_{\text{ref},t}[m,n,o]$ and $x_{\text{rx},t}[m,n,o]$, respectively. Hence, the TPSD planes are denoted by $\overline{\mathcal{S}}_{\text{ref},t}[h,k]$ and $\overline{\mathcal{S}}_{\text{rx},t}[h,k]$, respectively.

The local cross-correlations map of the anchor and distorted power spectra, $\zeta_t[h,k]$, is locally calculated within windows. The local cross-correlations plane is obtained as follows:

$$\zeta_t[h,k] = \frac{\sigma_{\overline{\mathcal{S}}_{\text{ref},t} \cdot \overline{\mathcal{S}}_{\text{rx},t}}[h,k] + C}{\sigma_{\overline{\mathcal{S}}_{\text{ref},t}}[h,k] \cdot \sigma_{\overline{\mathcal{S}}_{\text{rx},t}}[h,k] + C}, \quad (3)$$

where

$$\sigma_{\overline{\mathcal{S}}_{\text{X},t}}[h,k]$$
$$= \sqrt{\sum_{u=-d}^{d} \sum_{v=-d}^{d} \omega_{u,v}(\overline{\mathcal{S}}_{\text{X},t}[h+u,k+v] - \mu_{\overline{\mathcal{S}}_{\text{X},t}}[h,k])^2}, \quad (4)$$

$$\sigma_{\overline{\mathcal{S}}_{\text{X},t} \cdot \overline{\mathcal{S}}_{\text{Y},t}}[h,k]$$
$$= \sum_{u=-d}^{d} \sum_{v=-d}^{d} \omega_{u,v}(\overline{\mathcal{S}}_{\text{X},t}[h+u,k+v] - \mu_{\overline{\mathcal{S}}_{\text{X},t}}[h,k])$$
$$\times (\overline{\mathcal{S}}_{\text{Y},t}[h+u,k+v] - \mu_{\overline{\mathcal{S}}_{\text{Y},t}}[h,k]), \quad (5)$$

and

$$\mu_{\overline{\mathcal{S}}_{\text{X},t}}[h,k] = \sum_{u=-d}^{d} \sum_{v=-d}^{d} \omega_{u,v} \overline{\mathcal{S}}_{\text{X},t}[h+u,k+v]. \quad (6)$$

$\sigma_{\overline{\mathcal{S}}_{\text{ref},t} \cdot \overline{\mathcal{S}}_{\text{rx},t}}$ is the cross-covariance, $\mu_{\overline{\mathcal{S}}_{\text{ref},t}}$ and $\mu_{\overline{\mathcal{S}}_{\text{rx},t}}$ are the means, $\sigma_{\overline{\mathcal{S}}_{\text{ref},t}}$ and $\sigma_{\overline{\mathcal{S}}_{\text{rx},t}}$ are the standard deviations of $\overline{\mathcal{S}}_{\text{ref},t}$ and $\overline{\mathcal{S}}_{\text{rx},t}$, respectively, and $\omega$ is derived from 2D circular symmetric Gaussian weighting function with the window size of $11 \times 11$ ($d = 5$).

The term $\zeta_t[h,k]$ in (1) defines a 2D tempospatial full-reference perceptual quality map for tensor $t$ of the distorted video at every discrete frequency. In our implementation, 30 frames are grouped to form one tensor ($M = 1280, N = 720, O = 30$) and $C = 4.5 \times 10^{-4}$ is set to prevent instability when denominator is very close zero. The local cross-

**Fig. 3**. The incremental change in tempospatial PSD plane for the same video and same set of frames subject to different distortion levels. This example was taken from the Mobile LIVE database, sequence `Panning Under Oak`, frames $225 - 254$. The distortion magnitudes in the videos are as follows: $r1 > r2 > r3 > r4 > $ `Org` where `Org` is the anchor video free of distortion.

correlation map, $\zeta_t [h, k]$, is then averaged to obtain the tensor's perceptual quality score, $\overline{\zeta_t}$, as follows:

$$\overline{\zeta_t} = \frac{1}{MN} \sum_{\forall h} \sum_{\forall k} \zeta_t [h, k]. \tag{7}$$

For the temporal pooling of tensors to obtain the overall video quality score, we tested various pooling and statistical processing strategies. Mean pooling was chosen after it has empirically proven its superiority to other functions. Therefore, the overall video quality score is calculated by the average temporal quality of its tensors. That is,

$$\mathcal{P} = \left[ \operatorname*{E}_{\forall t} \left[ \overline{\zeta_t} \right] \right]^{\beta}, \tag{8}$$

where $\beta$ is an empirically determined sequence-dependent parameter.

Fig. 2 shows two simple examples to explain our proposed metric. Both sequences are composed of two frames only. The first frame, Frame 0, is identical in both examples. In Fig. 2a, the second frame, Frame 1, is identical to the first one. Only the distorted version is shown in the figure. In Fig. 2b, the edge in Frame 1 is shifted downwards to introduce a simple motion from the previous frame. The distorted version is shown in the figure. We show for each sequence the TPSD, $\overline{\mathcal{S}}_t [h, k]$, for both the distorted and anchor sequences. We also show the difference map between the two as well as local cross-correlations map, $\zeta_t [h, k]$. Moreover, Figure 3 shows the incremental change in the TPSD planes for different levels of distorted tensors with the same contents.

$\zeta_t [h, k]$ is a local cross-correlations map, which does not evaluate fidelity, it rather examines the contents in a certain frequency and quantifies the cross-correlation or consistency of contents in that frequency neighborhood with original contents. All the temporal and spatial contents corresponding to a certain frequency are unified within this 2D map. Every frequency spectrum in the original contents emits a certain optical energy to stimulate the HVS. A visual distortion will alter this energy in a certain way depending on the nature and severity of the distortion. This in turn causes discomfort and annoyance to viewers. In the context of visual masking, this framework models the visual sensitivity to distortions by estimating the power spectral cross-correlation, where at every frequency this local cross-correlation estimates the human visual discomfort in that frequency neighborhood. Quantifying the cross-correlation of spectral data in every frequency neighborhood measures the masking effect of the original contents (mask) in the presence of distortion (target). In other words, the local cross-correlation acts as a measure of annoyance or discomfort due to disruption of the original power spectrum caused by induced distortion. A high positive correlation indicates the contents to be similar, which yields little to no distortion to the viewer. Low positive and negative local cross-correlation values indicate a degradation in perceptual quality due to distortion. By averaging the map to obtain $\overline{\zeta_t}$, we incorporate the contribution to discomfort from every frequency. This averaging operation penalizes frequency spectra with low positive and negative cross-correlations by reducing the overall average for the whole tensor's perceptual quality.

## 3. EXPERIMENTS AND RESULTS

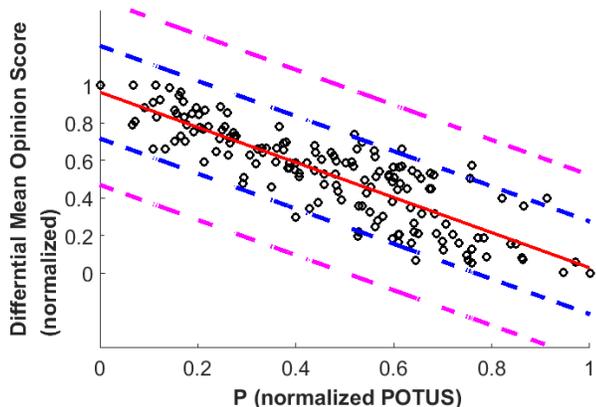We validate our proposed metric on the LIVE Mobile Video Quality Assessment database [14]. This database consists of 10 reference videos and 200 distorted videos. All videos are provided in `YUV420` format. They are 10 seconds in duration at a frame rate of 30 fps and a resolution of $1280 \times 720$.

Every anchor video in the databases has 20 different distorted versions as follows: four different levels of H.264 com-

pression artifacts, four different wireless packet loss levels, three different rate-adaptation patterns, five different temporal dynamics patterns, and 4 frame-freeze patterns. Compression artifacts were generated by using the JM reference implementation of the H.264 scalable video codec (SVC). Wireless packet loss patterns were simulated by transmitting H.264 SVC compressed video through a simulated wireless channel. Rate adaptation videos have a single rate switch in the video stream. Temporal dynamics videos contain multiple rate switches to test the effect of changes in video quality on the perceived quality. In addition, for the temporal dynamics distortion patterns in the database, we used only the last 210 frames to calculate overall video quality scores. This was motivated by the fact that DMOS values are mostly affected by the last a few seconds of the video [15]. We validated this choice by experimentally verifying that the correlation scores are maximum for all metrics using this range of frames after testing for all other combinations including the full set of frames.

The performance of our proposed metric is evaluated by calculating correlation scores with the DMOS scores reported in the database. Moreover, we compared our metric with commonly used and state of the art full-reference VQA metrics such as MOVIE [6], VQM [16], MS-SSIM [17], VIF [18], VSNR [19] and NQM [20]. PSNR is included as a baseline VQA metric.

In addition to the correlation score, we benchmark the computation time to calculate overall video quality score. We chose NQM and VIF to compare the computational time with our proposed method since these two metrics show a comparable performance to the proposed metric in terms of the correlation scores reported in Tables 1-2. All simulations were performed on a Windows PC with Intel Core i7-6700K CPU @ 4.00GHz, 32.0 GB RAM and MATLAB R2015(b).



**Fig. 4**. The predicted quality scores from the metric proposed in this work versus the reported DMOS for the all the sequences in the database. The blue and pink lines are P $\pm\sigma$ and P $\pm2\sigma$, respectively, where $\sigma$ is the data standard deviation.

**Table 1**. Pearson correlation coefficients (PCC) with the DMOS. (Co: Compression, Wl: Wireless channel packet loss, Ra: Rate adaptation, Td: Temporal dynamics)

| Distortion | Pearson Correlation Coefficients | | | | |
| --- | --- | --- | --- | --- | --- |
| | Co | Wl | Ra | Td | All |
| PSNR | 0.784 | 0.762 | 0.536 | 0.417 | 0.691 |
| VQM | 0.782 | 0.791 | 0.591 | 0.407 | 0.702 |
| MOVIE | 0.810 | 0.727 | 0.681 | 0.244 | 0.716 |
| MS-SSIM | 0.766 | 0.771 | 0.709 | 0.407 | 0.708 |
| VIF | 0.883 | 0.898 | 0.664 | 0.105 | 0.787 |
| VSNR | 0.849 | 0.849 | 0.658 | 0.427 | 0.759 |
| NQM | 0.832 | 0.874 | 0.677 | 0.365 | 0.762 |
| Proposed | **0.951** | **0.949** | **0.856** | **0.800** | **0.850** |

**Table 2**. Spearman correlation coefficients (SCC) with the DMOS. (Co: Compression, Wl: Wireless channel packet loss, Ra: Rate adaptation, Td: Temporal dynamics)

| Distortion | Spearman Correlation Coefficients | | | | |
| --- | --- | --- | --- | --- | --- |
| | Co | Wl | Ra | Td | All |
| PSNR | 0.819 | 0.793 | 0.598 | 0.372 | 0.678 |
| VQM | 0.772 | 0.776 | 0.648 | 0.386 | 0.695 |
| MOVIE | 0.774 | 0.651 | 0.720 | 0.158 | 0.642 |
| MS-SSIM | 0.804 | 0.813 | 0.738 | 0.397 | 0.743 |
| VIF | 0.861 | 0.874 | 0.639 | 0.124 | 0.744 |
| VSNR | 0.874 | 0.856 | 0.674 | 0.317 | 0.752 |
| NQM | 0.850 | 0.899 | 0.678 | 0.238 | 0.749 |
| Proposed | **0.959** | **0.952** | **0.879** | **0.811** | **0.858** |

### 3.1. Results and Analysis

Figure 4 shows the scatter plot of predicted overall video quality score from the proposed method versus DMOS reported in the database. Most of the scatter points are located within one standard deviation boundaries (blue lines). Outliers are also very close to the $\mathcal{P} \pm \sigma$ lines. This shows that our predicted quality scores are highly correlated with the subjective human scores of video quality.

Table 1-2 show Pearson's Correlation Coefficients (PCC) and Spearman's Correlation Coefficients (SCC) calculated between predicted quality scores and DMOS in the database. We report PCC and SCC values for each distortion type and with all the videos in the database. The bolded values represent the highest value in each column. For both PCC and SCC, our proposed method outperforms all other VQA metrics by a significant margin for the whole database as well as for every distortion type. In particular, our metric outperforms the second best in compression artifacts (VSNR) by 0.085 in terms of SCC. For wireless packet loss distortions, the proposed metric outperforms the second best (NQM) by 0.053 in terms of SCC. Furthermore, PCC and SCC values from other VQA metrics significantly drop when algorithms are applied

**Table 3**. Computation time to calculate the quality score of 120 frames of `Harmonicat` video (frames $201 - 320$).

| Time [sec] | Computation time | | |
|---|---|---|---|
| | VIF | NQM | Proposed |
| | 255.729 | 59.490 | 15.030 |

to rate adaptation and temporal dynamics distortions. However, our proposed metric shows a robust performance on both distortion types in terms of PCC and SCC. It shows above 0.8 values of PCC, SCC on rate adaptation and temporal dynamics. Since this framework includes both spatial and temporal features via 3D PSD processing, the algorithm effectively captures the impact of dynamic rate changes to human perceived video quality.

Table 3 shows the computational efficiency of our approach compared to other full-reference VQA metrics. This metric only needs 5.88% of computational time required by VIF and 25.26% of computational time required by NQM. DFT is one of the simplest domain transform operations and our framework processes a 2D time-aggregated PSD plane for a tensor of frames instead of individual frame processing, which decreases computational burden.

## 4. CONCLUSION

We propose a full-reference PVQA metric through 3D PSD analysis. In particular, we utilize 2D time-aggregated PSD plane to obtain tempospatial power features and calculate cross-correlation with the reference to quantify the effect of distortion on perceived video quality. We thoroughly evaluate the performance of the proposed metric in terms of correlation with human mean opinion scores of video quality. The results show competitive correlations compared with the state of the art techniques. This work does not make any assumption on coding conditions or video sequence. Furthermore, Our proposed metric has a low computational complexity, which makes it feasible for real-time application. We believe that this work to be a significant step towards understanding the relationship between PSD and perceived quality.

## 5. REFERENCES

[1] "The zettabyte era: Trends and analysis," http://www.cisco.com/c/en/us/solutions/collateral/service-provider/visual-networking-index-vni/vni-hyperconnectivity-wp.pdf, Feb 2016, Accessed: 2014-12-01.

[2] "Alliance for open media - aom in the news," http://aomedia.org/aom-in-the-news/, (Accessed on 10/20/2016).

[3] "The state of video codecs 2016 - streaming media magazine," http://www.streamingmedia.com/Articles/Editorial/Featured-Articles/The-State-of-Video-Codecs-2016-110117.aspx, 2016, (Accessed on 10/29/2016).

[4] "Audio video coding standard workgroup of china," http://www.avs.org.cn/english/Achievement.asp#2.1, (Accessed on 10/20/2016).

[5] "How consumers judge their viewing experience," http://lp.conviva.com/rs/901-ZND-194/images/CSR2015_HowConsumersJudgeTheirViewingExperience_Final.pdf, 2015, (Accessed on 11/28/2016).

[6] K. Seshadrinathan and A.C. Bovik, "Motion tuned spatio-temporal quality assessment of natural videos," *Image Processing, IEEE Transactions on*, vol. 19, no. 2, pp. 335–350, Feb 2010.

[7] M. A. Aabed and G. AlRegib, "Reduced-reference perceptual quality assessment for video streaming," in *International Conference on Image Processing (ICIP), 2015 IEEE*, Sept 2015.

[8] 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, April 2013.

[9] 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, March 2014.

[10] Y. Li, L. M. Po, C. H. Cheung, X. Xu, L. Feng, F. Yuan, and K. W. Cheung, "No-reference video quality assessment with 3d shearlet transform and convolutional neural networks," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 26, no. 6, pp. 1044–1057, June 2016.

[11] Farah Torkamani-Azar, Hassan Imani, and Hossein Fathollahian, "Video quality measurement based on 3-d. singular value decomposition," *Journal of Visual Communication and Image Representation*, vol. 27, pp. 1 – 6, 2015.

[12] A. Mittal, M. A. Saad, and A. C. Bovik, "A completely blind video integrity oracle," *IEEE Transactions on Image Processing*, vol. 25, no. 1, pp. 289–300, Jan 2016.

[13] Antonio Torralba and Aude Oliva, "Statistics of natural image categories," *Network: Computation in Neural Systems*, vol. 14, no. 3, pp. 391–412, 2003, PMID: 12938764.

[14] A. K. Moorthy, L. K. Choi, A. C. Bovik, and G. de Veciana, "Video quality assessment on mobile devices: Subjective, behavioral and objective studies," *IEEE Journal of Selected Topics in Signal Processing*, vol. 6, no. 6, pp. 652–671, Oct 2012.

[15] Iain E. Richardson, *Video Codec Design: Developing Image and Video Compression Systems*, John Wiley & Sons, Inc., New York, NY, USA, 2002.

[16] M. H. Pinson and S. Wolf, "A new standardized method for objectively measuring video quality," *IEEE Transactions on Broadcasting*, vol. 50, no. 3, pp. 312–322, Sept 2004.

[17] Z. Wang, E. P. Simoncelli, and A. C. Bovik, "Multiscale structural similarity for image quality assessment," in *The Thrity-Seventh Asilomar Conference on Signals, Systems Computers, 2003*, Nov 2003, pp. 398–1402 Vol.2.

[18] H. R. Sheikh and A. C. Bovik, "Image information and visual quality," *IEEE Transactions on Image Processing*, vol. 15, no. 2, pp. 430–444, Feb 2006.

[19] D. M. Chandler and S. S. Hemami, "Vsnr: A wavelet-based visual signal-to-noise ratio for natural images," *IEEE Transactions on Image Processing*, vol. 16, pp. 2284–2298, Sept 2007.

[20] N. Damera-Venkata, T. D. Kite, W. S. Geisler, B. L. Evans, and A. C. Bovik, "Image quality assessment based on a degradation model," *IEEE Transactions on Image Processing*, vol. 9, no. 4, pp. 636–650, Apr 2000.