

Subjective and Objective Study of the Relation Between 3D and 2D Views Based on Depth and Bitrate

Balasubramanyam Appina, Manasa K, and Sumohana S. Channappayya
*Lab for Video and Image Analysis (LFOVIA), Department of Electrical Engineering,
Indian Institute of Technology Hyderabad, India, 502285.*
e-mail: {ee13m14p100001, ee12p1002, sumohana}@iith.ac.in.

Abstract

The tremendous growth in 3D (stereo) imaging and display technologies has led to stereoscopic content (video and image) becoming increasingly popular. However, both the subjective and the objective evaluation of stereoscopic video content has not kept pace with the rapid growth of the content. Further, the availability of standard stereoscopic video databases is also quite limited. In this work, we attempt to alleviate these shortcomings. We present a stereoscopic video database and its subjective evaluation. We have created a database containing a set of 144 distorted videos. We limit our attention to H.264 compression artifacts. The distorted videos were generated using 6 uncompressed pristine videos of left and right views originally created by Ecole Polytechnique Federal De Lausanne (EPFL)[1]. The reference video sequences contain a good combination of texture, motion, depth information and we divided these videos into 2 groups based on depth information. Further, 19 subjects participated in the subjective assessment task. Based on the subjective study, we have formulated a conditional relation between the 2D and stereoscopic subjective scores as a function of compression rate and depth range. We have also evaluated the performance of popular 2D and 3D image/video quality assessment (I/VQA) algorithms on our database.

Introduction

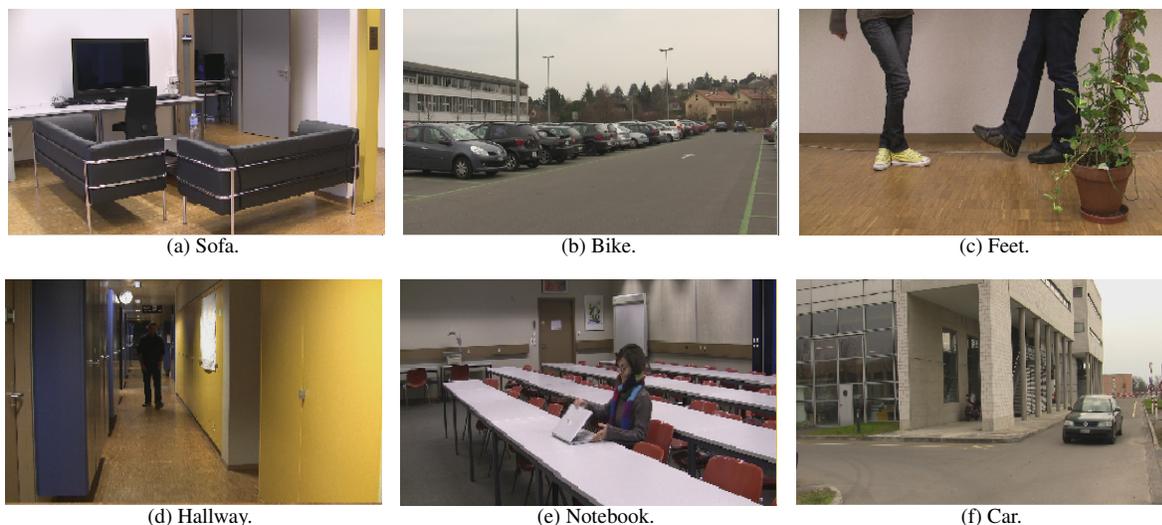
With the rapid advancements in 3D video technology, the industry and consumer experiences are improving in a tremendous way. According to the recent survey by Motion Picture Association of America (MPAA) [2], the US revenue from 3D film industry has risen to 16% in 2013 and one third of movie lovers watch at least one 3D movie in a month. The primary reason for this incredible increase could be attributed to the depth-enhanced viewing experience. This has led to the movie and gaming industries investing a significant amount of resources on the creation of 3D content.

The creation of multimedia content happens over several processing stages (such as sampling, quantization, demosaicing etc.), each of which could potentially degrade the perceptual quality of the content. Compression artifacts are a very common cause of quality degradation. In this work, we focus our attention only on compression artifacts. Given that most of this content is meant for human consumption, the most relevant and consistent method to evaluate video quality is via subjective assessment. While subjective assessment is cumbersome, expensive and time consuming, this data is very essential to test the performance of objective VQA algorithms.

In this paper, we present the subjective quality assessment of stereoscopic videos. The subjective assessment for stereoscopic videos is different from that of the 2D video, as the stereoscopic video consists of two views: left view and right view. These two views contribute to the perception of depth. Therefore, the overall quality of a stereoscopic video is a function of the individual qualities of the constituent left and right views.

Goldmann et al. [1] created a database to study the effect of variation in the distance between the camera and the objects on perception. Ha et al. [3] performed a subjective study on a stereoscopic video data set based on the consideration of visual quality, depth perception, visual comfort and overall quality. Their work mainly focused on the perception of depth information without considering the distortion in the videos. Hewage et al. [4] conducted a subjective study to explore the effect of random packet loss artifacts on the overall perceptual quality of stereoscopic video. Aflaki et al. [5] have performed a subjective study to explore the effects of asymmetric encoding (different rates and resolutions assigned to the left and right views) of a stereoscopic video. They conclude that asymmetric encoding offers bitrate savings compared to the symmetric case. They also concluded that PSNR is not a good objective measure for analyzing blocking artifacts and blurriness. Urvoy et al. [6] created a symmetrically distorted stereoscopic video dataset composed of H.264, JPEG 2000 as compression artifacts. However, this database does not consider the important case of asymmetric distortion of the stereoscopic views. While these subjective studies have considered either symmetric or asymmetric distortions, they have not considered the relationship between stereoscopic views and 2D views as a function of compression rate and depth. De Silva et al. [7] proposed a FR 3D VQA based on measuring the structural distortion, blur strength and content complexity. The structural distortion strength is computed by calculating the similarity measurement between reference and distorted frames and further, disturbances in edge strength is computed to measure the blur strength. The content complexity is measured by calculating the spatial index (SI) and temporal index (TI) based on ITU recommendation P.910 of a 3D view. Cheng et al. [8] created a publicly available uncompressed stereoscopic dataset. Video sequences have a resolution of 1920×1080 and a frame rate of 25 fps. Chen et al. [9] created a H.264 stereoscopic dataset to explore the stereoscopic quality effect parameters. They utilized the EPFL stereoscopic video sequences to perform the study and video are resized to 720×480 . They concluded that subjective scores have unique trend in spatial quality but it has a relatively different trend in depth quality scores.

Figure 1: One frame from each pristine video.



While our work is similar in philosophy to [5], [6], [7] we would like to highlight our contributions: a) creation of a stereoscopic video database that would be made freely available to the research community, b) a study of the relationship between the stereoscopic subjective scores and 2D subjective scores, c) a performance evaluation of popular 2D and 3D I/VQA algorithms on our database, and d) an exploration of the effect of depth and compression rate on the perceptual quality of stereoscopic video.

Database Description

In this section we describe the generation of the video sequences used in our study, starting with a description of the pristine or reference sequences.

Pristine Sequences

Goldmann et al. at EPFL [1] created an open source database to highlight the effect of viewing distance variation between the camera and objects. We used the same reference videos as those in the EPFL study. There are six pristine videos per view (left view, right view) in the database. Fig. 1 shows one frame of each reference sequence in the database.

This database consists of a collection of indoor and outdoor scenes with varying range of color, texture and objects. These videos are captured with identical camcorders placed horizontally with the separation continuously adjustable in the range 7–50 cm. The videos have resolutions varying from 1836×1056 to 1900×1054 pixels and a frame rate of 25 fps. Each video is 10 seconds in duration and is placed in an *avi* container. The camcorders were controlled by a remote to account for any temporal mismatch.

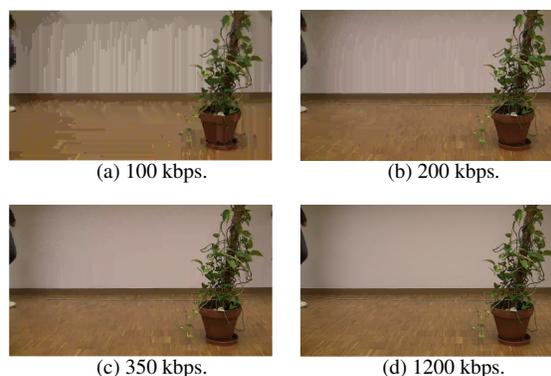
We grouped the 6 videos into two categories based on the depth content in them. *Sofa*, *Feet*, *Hallway*, *Notebook* sequences form group I having lower depth (3m - 10m). The *Bike* and *Car* sequences fall into group II having higher depth range (> 100m).

Test Sequences

The pristine sequences that were in the *avi* format were converted to the YUV 4:2:0 format using the open-source *ffmpeg* ap-

plication [10].

Figure 2: One frame from the *Feet* sequence for different compression rates.



As mentioned previously, the pristine videos had different resolutions, of which a majority were at a resolution of 1836×1056 pixels. To maintain consistency, videos at other resolutions were resized to 1836×1056 pixels using *ffmpeg*.

We generated 4 test sequences from each reference video using H.264 compression. We used a variety of compression rates (100 kbps, 200 kbps, 350 kbps and 1200 kbps) to cover a wide range of possible video transmission link rates. We capped our rate at 1200 kbps because the perceptual quality variation was not significant beyond this rate. Fig. 2 shows the frames of the pristine *Feet* sequence encoded at varying compression rates. The compression was done by the *ffmpeg* software using *libx264* at the following settings: GOP length of 250 frames (default), CABAC encoder, flags and loop filter enabled. The compression rate was fixed using the *maxrate* parameter.

Overall, there are 24 test sequences and 6 reference sequences per view. The left and right views were combined to form symmetric and asymmetric sets. We would like to recall that in the symmetric set, the left and right views have been encoded at the same bitrate. In the asymmetric set, the left and right views of

a video are encoded at different bitrates. The symmetric set contains 24 videos (6 reference videos encoded at 4 different rates). The asymmetric set has 120 videos (out of the 150 possible permutations, 24 belong to the symmetric set, 6 are reference pairs, and the remaining fall into the asymmetric set).

Subjective Study

Display Settings

We used a Samsung display of 32 inches (81.28 cm) with a screen resolution of 1366×768 pixels for our subjective study. The distance between the observer and screen was fixed at 1.5 meters which is 3 times the height of screen and the observer was seated at a height of 20.5 inches (52 cm) as shown in Fig. 3. The rest of the settings adhered to the ITU-R recommendations for subjective quality evaluation [11]. The stereoscopic videos were played using a NVIDIA stereoscopic player [12].

Assessment Method

We used the Single Stimulus Continuous Quality Evaluation (SSCQE) method to obtain the subjective rating of the videos. Our subjective study involved 19 subjects, gender distribution was not limited in our study and the average age of all observers is 24 years.



Figure 3: Subjective equipment setup.

A demo sequence that is representative of the quality variability in the distorted videos was first shown to the subjects. The subjective analysis was conducted in two sessions of 30 minutes each. During the first session the subjects were shown the left and right views of the 2D video. The videos were arranged in a random order of varying compression rates, and it was ensured that there were no repetition of video sequences. In the second session, the subjects were trained to perceive the stereoscopic content and asked to rate the stereoscopic videos. For stereoscopic quality evaluation, the subjects wore a pair of anaglyph glasses and the stereo videos were rendered using a NVIDIA stereoscopic player.

The subjective rating given to the video is according to the ITU-R ACR scale, which ranges from 1 - 5 (1 - bad, 2 - poor, 3 - fair, 4 - good, 5 - Excellent). Non-integer ratings were also allowed.

Subjective Scores Analysis

Subjective data handling

To process the subjective scores we followed the ITU-R recommendations [11][13]. We have 150 (symmetric + asymmetric) scores for a stereoscopic video set and 30 scores for each 2D view (left view and right view). First, we compute difference scores between the test video and reference video. These scores are computed by subtracting the quality score assigned by the subject to a test video from the quality score assigned by the same subject to

the corresponding reference video.

$$d_{ij} = s_{ij_{ref}} - s_{ij}, \quad (1)$$

where i indicates the subject and j indicates the video sequence id. The difference scores for the reference videos are not considered for analysis. The z -scores are computed by calculating the mean (μ_i) and standard deviations (σ_i) from difference scores for each subject. The z_{ij} scores are given by

$$\mu_i = \frac{\sum_{j=1}^{N_j} d_{ij}}{N_j}, \quad (2)$$

$$\sigma_i = \sqrt{\frac{\sum_{j=1}^{N_j} (d_{ij} - \mu_i)^2}{N_j - 1}}, \quad (3)$$

$$z_{ij} = \frac{d_{ij} - \mu_i}{\sigma_i}, \quad (4)$$

where N_j is the number of videos rated by the subject i . For the stereoscopic case, $N_j = 150$ and for the 2D cases, $N_j = 30$ for each view.

To remove outliers we followed the ITU-R BT 500.11 recommendations for observer screening. Observers are discarded if they exhibit a strong shift of votes compared to the average behaviour.

In our analysis no outliers were found.

The z -scores lie in the range of $[-3,3]$ which was scaled to $[0,100]$ by

$$z'_{ij} = \frac{100(z_{ij} + 3)}{6}, \quad (5)$$

The final step in subjective processing is calculation of DMOS scores. DMOS is calculated by taking the mean of the rescaled z -scores across all the subjects per video.

$$DMOS_j = \frac{\sum_{i=1}^M z'_{ij}}{M}, \quad (6)$$

where $M = 19$. The range of DMOS values obtained for stereoscopic set is $[79.73 \ 28.9]$. Similarly, for the left video set the range is $[73.99 \ 26.47]$ while it is $[72.36 \ 28.65]$ for the right video set.

Performance Evaluation

Subjective Score based Evaluation

Let L_j, R_j be the DMOS for the left and right views respectively for a video j . The average of the left and right view DMOS, V_j is given by

$$V_j = \frac{L_j + R_j}{2}. \quad (7)$$

Table 1 shows the correlation between V_j of a video with the corresponding stereoscopic DMOS. As defined earlier, in the symmetric case both views having same compression rate while the asymmetric case stands for different compression rates in the left and right view. For instance, the asymmetric case for 100 kbps compression rate denotes the compression rate of one of the views being fixed at 100 kbps and the other view's rate being varied for all combinations and vice versa.

Table 1: Correlation between the the average DMOS of left and right views V_j and the stereoscopic DMOS across varying compression rates

Compression rates	Symmetric	Asymmetric
100 kbps	0.280	0.563
200 kbps	0.939	0.912
350 kbps	0.872	0.934
1200 kbps	0.081	0.875

Table 2: Correlation scores of average DMOS, V_j and the stereoscopic DMOS as a function of depth and compression rates for videos belonging to group I and group II.

Compression rates	Asymmetric	
	I	II
100 kbps	0.48	0.62
200 kbps	0.77	0.92
350 kbps	0.82	0.96
1200 kbps	0.77	0.94

From Table 1 it is clearly seen that V_j is not a representative of the stereoscopic quality across the compression rates. Table 2 shows the correlation values between V_j and the stereoscopic DMOS of group I (lower depth range) and group II (higher depth range) for different compression rates.

Table 3: Comparison of stereoscopic DMOS for asymmetric video set of group I videos (lower depth range).

Compression rates	Asymmetric			
	100 kbps	200 kbps	350 kbps	1200 kbps
100 kbps		65.06	66.16	59.06
200 kbps	70.87		54.45	48.29
350 kbps	66.88	52.18		36.29
1200 kbps	62.05	49.25	42.65	

We present the following hypothesis to explain the performance of V_j as a stereoscopic quality metric. When a video is visually very annoying the viewer gets accustomed to it and tries to extract information from the given quality video. In case of stereoscopic video this information can be the depth range. As V_j does not have the effect of depth in it (since it the average of 2D scores), the correlation is low for higher compression rates. However, in the case of medium compression rate, the scenes are neither visually too annoying nor very good and hence results in viewer dilemma. Owing to this constraint, the viewer does not attempt to infer additional information at this rate. Therefore, as the compression rate decreases the correlation increases. When the video is of very high quality (lower compression rate), the viewer tries to capture additional information from the scene which is again depth in the stereoscopic case, and hence the correlation decreases.

Tables 3 and 4 show the average stereoscopic DMOS values for the asymmetric video sets for groups I and II respectively. It is clear that the stereoscopic DMOS for the group I videos are high compared to the group II videos. In group II videos the depth range is high which results in lower viewing precision of the objects in the scene. Therefore, the distortions at higher depth range are not easily perceived. Thus, for a given compression rate, we can conclude that the DMOS for stereoscopic videos with higher depth range is always less than the videos with lower depth range.

Table 4: Comparison of stereoscopic DMOS for asymmetric sets of group II videos (higher depth range).

Compression rates	Asymmetric			
	100 kbps	200 kbps	350 kbps	1200 kbps
100 kbps		56.89	57.31	54.04
200 kbps	63.27		45.01	45.03
350 kbps	61.71	44.22		35.78
1200 kbps	59.88	43.16	36.5	

Table 5: Comparison of the performance of the 2D metrics on the left and right views of the database across the compression rates.

Algorithm	Compression rates			
	100 kbps	200 kbps	350 kbps	1200 kbps
PSNR [14]	0.64	0.53	0.29	0.76
VSNR [15]	0.19	0.77	0.11	0.12
SSIM [16]	0.69	0.73	0.45	0.77
FSIM [17]	0.69	0.78	0.49	0.76
STMAD [18]	0.27	0.42	0.62	0.72
BVQM [20]	0.74	0.93	0.89	0.76

Objective Score based Evaluation

In order to test the efficacy of popular 2D and 3D objective I/VQA metrics on stereoscopic video, they were evaluated on the stereoscopic database we have created. Standard measures of performance such as Spearman Rank Order Correlation Coefficient (SROCC) and Linear Correlation Coefficient (LCC) were used. Table 5 shows the performance of the 2D I/VQA metrics on the left and right view videos of the database. A non-linear regression on the VQA scores is done using the logistic function mentioned in [21] and LCC is computed between the fitted objective scores and the DMOS. PSNR [14], VSNR [15], SSIM [16], FSIM [17] are image metrics and they are applied on a frame by frame basis and averaged. ST-MAD [18] and BVQM [19, 20] are 2D video metrics. The Chen et al. [22] and STRIQE [23] are 3D IQA metrics. Due to the high computational time complexity the performance of the publicly available 3D NR IQA metrics [24], [25] were not tested on our dataset.. Table 6 illustrates the performance of the 2D and 3D I/VQA metrics across compression rates for stereoscopic videos in the database.

Conclusions and Future work

The purpose of this study was to create a H.264 compressed stereoscopic video dataset. The created stereoscopic video database composed of 144 videos was created using the 6 pristine videos from the EPFL database [1], and compressed at 4 compression rates. The subjective study was done on these videos by 19 subjects. We tested the efficacy of several 2D I/VQA algorithms on the proposed database.

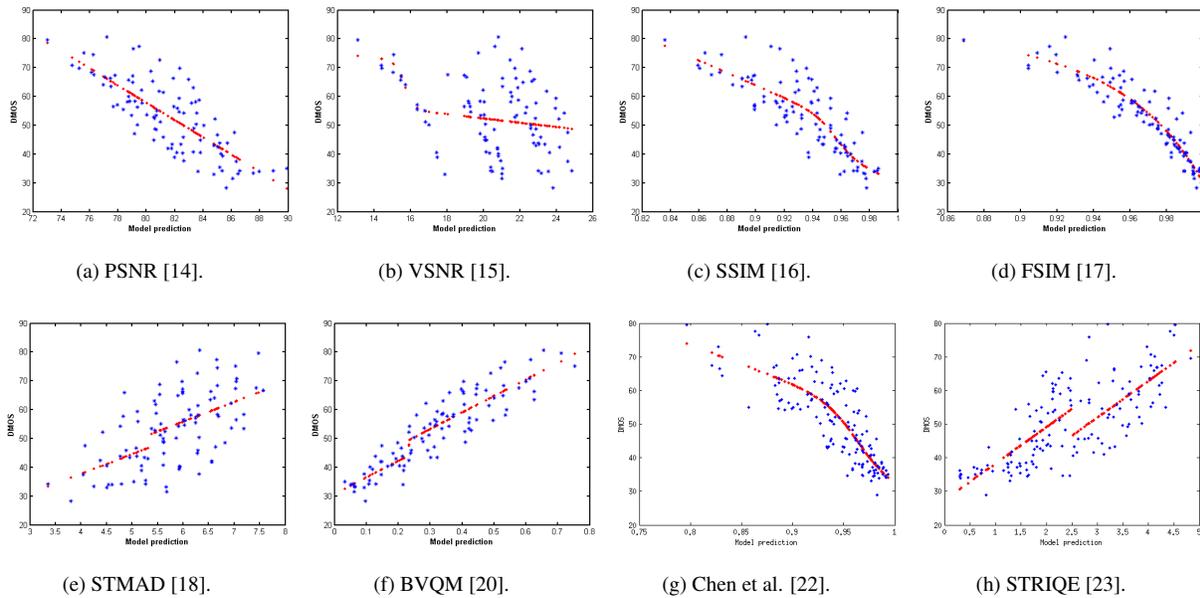
From the analysis of the subjective scores we made the following conclusions: i) the average DMOS from the left and right views V_j is not a representative of the stereoscopic DMOS, ii) depth plays a role at very high and very low compression rates. Therefore, the 2D and stereoscopic I/VQA algorithms do not perform well at high and low compression rates, iii) the study on the correlation values has depicted that at a given compression rate, the videos with higher depth range have better visual quality compared to that of lower depth range ones. The objective VQA perform better on the videos having higher depth range.

We plan to make the database and the DMOS values avail-

Table 6: Comparison of the performance of the average of the 2D objective metrics on the left and right views of the stereoscopic database across the compression rates for group I (lower depth range) and group II (higher depth range videos) - Linear Correlation Coefficient.

Algorithm	Asymmetric							
	Compression rates							
	100 kbps		200 kbps		350 kbps		1200 kbps	
	I	II	I	II	I	II	I	II
PSNR [14]	0.71	0.55	0.63	0.74	0.70	0.37	0.68	0.64
VSNR [15]	0.66	0.45	0.54	0.71	0.57	0.68	0.53	0.50
SSIM [16]	0.71	0.86	0.77	0.90	0.85	0.96	0.81	0.95
FSIM [17]	0.97	0.94	0.95	0.96	0.92	0.92	0.98	0.98
BVQM [19, 20]	0.73	0.74	0.88	0.92	0.92	0.96	0.94	0.97
STMAD [18]	0.66	0.69	0.75	0.89	0.77	0.90	0.80	0.82
<i>Chen et al. [22]</i>	0.62	0.92	0.64	0.98	0.77	0.98	0.88	0.97
<i>STRIQUE [23]</i>	0.65	0.77	0.71	0.89	0.77	0.91	0.74	0.87

Figure 4: Scatter plots of objective scores versus the stereo DMOS over all the videos in the database.



able publicly to the research community. A sample set of videos can be found at *LFOVIA's* home page [26].

References

- [1] Goldmann L, De Simone F, Ebrahimi T. A comprehensive database and subjective evaluation methodology for quality of experience in stereoscopic video. IS&T, SPIE, 2010, pg. 75260S.
- [2] Motion Picture Association of America. Theatrical market statistics 2013.
- [3] K Ha and M Kim. A perceptual quality assessment metric using temporal complexity and disparity information for stereoscopic video. ICIP 2011 IEEE, pg. 25252528.
- [4] C Hewage, M Martini, M Brandas, and D De Silva. A study on the perceived quality of 3d video subject to packet losses. ICC 2013 IEEE, pg. 662666.
- [5] P Aflaki, M M Hannuksela, J Hakkinen, P Lindroos, and M Gabbouj. Subjective study on compressed asymmetric stereoscopic video. ICIP 2010 IEEE, pg. 4021.
- [6] M Urvoy, M Barkowsky, R Cousseau, Y Koudota, V Ricorde, P Le Callet, J Gutierrez, and N Garcia, Nama3ds1-cospad1: Subjective video quality assessment database on coding conditions introducing freely available high quality 3d stereoscopic sequences. QoMEX 2012 IEEE, pg. 109.
- [7] V D Silva, H K Arachchi, E Ekmekcioglu, and A Kondoz. Toward an impairment metric for stereoscopic video: A full-reference video quality metric to assess compressed stereoscopic video. IEEE TIP 2013, pg. 3392.
- [8] Cheng E, Burton P, Burton J, Joseski A, Burnett I. RMIT3DV: Pre-announcement of a creative commons uncompressed HD 3D video database. QoMEX IEEE 2012, pg. 212.
- [9] M J Chen, D K Kwon, and A Bovik. Study of subject agreement on stereoscopic video quality. SSIAP IEEE 2012, pg. 173.
- [10] <https://www.ffmpeg.org/>
- [11] Int. telecommun. union, methodology for the subjective assessment of the quality of television pictures itu-r recommendation, pg. BT.50011, 2000.
- [12] <http://www.nvidia.com/object/3d-vision-video-player-1.7.5-driver.html>
- [13] K Seshadrinathan, R Soundararajan, A C Bovik, and L K Cormack. Study of subjective and objective quality assessment of video. IEEE

- TIP 2010, pg. 14271441.
- [14] H R Sheikh, M F Sabir, and A C Bovik. A statistical evaluation of recent full reference image quality assessment algorithms. IEEE TIP 2006, pg. 3440.
- [15] D M Chandler and S S Hemami. Vsnr: A wavelet-based visual signal-to-noise ratio for natural images. IEEE TIP 2007, pg. 2284.
- [16] Z Wang, A C Bovik, H R Sheikh, and E P Simoncelli. Image quality assessment: From error visibility to structural similarity. IEEE TIP 2004, pg. 600.
- [17] L Zhang, D Zhang, and X Mou. Fsim: a feature similarity index for image quality assessment. IEEE TIP, pg. 2378.
- [18] P V Vu, C T Vu, and D M Chandler. A spatiotemporal mostapparent-distortion model for video quality assessment. ICIP 2011 IEEE.
- [19] Vqm software: <http://www.its.bldrdoc.gov/n3/video/vqmsoftware.htm>.
- [20] M H Pinson and S Wolf. A new standardized method for objectively measuring video quality. Broadcasting IEEE 2004, pg. 312.
- [21] (2000) final report from the video quality experts group on the validation of objective quality metrics for video quality assessment.
- [22] M J Chen, C C Su, D K Kwon, L K Cormack, and A C Bovik. Full-reference quality assessment of stereopairs accounting for rivalry. SIP Image Communication 2016, pg. 1143.
- [23] S Khan Md, B Appina, and S Channappayya. Full-reference stereo image quality assessment using natural stereo scene statistics. SPL IEEE 2015, pg. 1985.
- [24] Appina B, Khan S, Channappayya S S. No-reference Stereoscopic Image Quality Assessment Using Natural Scene Statistics. SIP Image Communication 2016, pg. 43.
- [25] M J Chen, L K Cormack, A C Bovik. No-reference quality assessment of natural stereopairs. IEEE TIP 2013. pg. 3379.
- [26] <http://www.iith.ac.in/~Ifovia>.

Author Biography

- **Balasubramanyam Appina** is currently a Research Scholar in the Electrical Engineering department of Indian Institute of Technology, Hyderabad. He is doing research in 2D & 3D image and video quality assessment using natural scene statistical (NSS) features. Received gold medal recipient from International Federation of Inventors Associations for best National Innovation (2016).
- **Manasa K** is currently a Research Scholar in the Electrical Engineering department of Indian Institute of Technology, Hyderabad. Her research interests include image and video quality assessment, statistical modeling of videos to measure quality, modeling quality of experience and quality of service. She is a recipient of the Excellence in Research award at IIT Hyderabad (2015) and the Above and Beyond Award at PEC (2011). She was a R&D consultant at Sony India Software Center, India (2015-2016).
- **Sumohana S. Channappayya** is currently an Assistant Professor of Electrical Engineering with IIT Hyderabad, where he directs the Laboratory for Video and Image Analysis. He completed his Ph.D. degree in electrical and computer engineering from The University of Texas at Austin, in 2007. His research interests include image and video quality assessment, multimedia communication, and biomedical imaging. He was a recipient of the Excellence in Teaching Award at IIT Hyderabad (2013).