

PERCEPTUAL QUALITY METRIC FOR COMPRESSED VIDEOS

EePing Ong*, Xiaokang Yang#, Weisi Lin*, Zhongkang Lu*, and Susu Yao*

*Institute For Infocomm Research, 21 Heng Mui Keng Terrace, Singapore 119613

#Shanghai Jiao Tong University, 1954 Hua Shan Road, Shanghai 200030, China

Email: {epong, wslin, zklu, ssyao}@i2r.a-star.edu.sg; xkyang@cdtv.org.cn

ABSTRACT

This paper proposed an objective perceptual video quality metric to automatically assess the perceived quality of digital videos. Traditionally, Peak signal-to-noise ratio (PSNR) has been used to represent the quality of a compressed video sequence. However, PSNR has been found to correlate poorly with subjective quality ratings, particularly at much lower bit rates and frame rates. In this paper, computational models have been applied to emulate human visual perception based on block-fidelity, content richness fidelity, spatial-textural, colour, and temporal maskings. The proposed video quality metric has been tested on CIF and QCIF video sequences compressed at various bit rates and frame rates. It has been shown to give significantly better correlations to human perception than peak signal-to-noise ratio (PSNR).

1 INTRODUCTION

Besides on-line and off-line visual quality evaluation, how distortion is gauged also plays a determinative role in shaping most algorithms for image/video manipulations. Visual quality control within an encoder and distortion assessment for decoded signal is particularly of interests due to the widespread applications of H.26x/MPEG-x compression. Since human eyes are the end receiver of most decoded images/videos, it is desirable to develop visual quality metrics that correlate better with human's visual perception than the conventional PSNR measure (which has been found to correlate poorly with subjective quality ratings [3]).

For better prediction of quality for decoded visual signal, a number of approaches have been tried to model the temporal, spatial and masking characteristics of human vision [15, 16, 9, 14, 7], to evaluate common coding artefacts [17, 5], and also to combine these two paradigms [12, 18]. However, the evaluation so far has been concentrated in TV types of signal (e.g., [14]). Characteristics of human vision have also been modelled for image and video compression [2, 4]). There will be increasingly more applications and services of mobile visual signals, and therefore this creates the need for measuring/monitoring the quality of images/video coded at low bit rates. The lower end of bit rates for wireless multimedia may be in the region of 10 kb/s for QCIF images/videos. It is expected that the lower the coding bit rate, the worse the conventional pixel-wise error measures (e.g., MSE, PSNR) perform in terms of visual quality assessment. It has been shown during this work that the PSNR quality prediction accuracy in low bit rate situations is worse than that reported in [14].

In the rest of this paper, Section 2 presents a description of the proposed video quality metric, while Section 3 demonstrates the performance of the proposed scheme in low bit rate coded QCIF and CIF videos. The last section concludes this paper.

2 PROPOSED VIDEO QUALITY METRIC

2.1 Video Quality Metric

The overall objective video rating for a colour video sequence, Q , is given by a weighted averaging of the objective video quality rating for each colour's $q_j(t)$, for $j=1, \dots, n$, where n is the maximum number of colour components, and can be expressed as:

$$Q = \sum_{j=1}^n \alpha_j \left(\frac{\sum_{t=i(f_f/f_r)}^N [q_j(t)]}{N_t} \right), \quad i=1,2,\dots$$

where α_j denotes the weighting for each colour components, N is the total number of frames in the original video sequence, $N_t = N/(f_f/f_r)$ is the total number of frames in the coded video sequence, f_r is the frame rate at which the video is being coded, f_f is the full frame rate of the original video sequence, and $q_j(t)$ is the objective video quality rating for colour component of each frame:

$$q_j(t) = D_j(t) \cdot F_{BI,j}(t) \cdot F_{RF,j}(t)$$

where D is the distortion-invisibility (derived from spatial-textural, colour and temporal maskings), F_{BF} is the block-fidelity, and F_{RF} is the content richness fidelity. The latter two terms are global measures that modulate the final distortion-invisibility value to give the video quality measure for each frame. These global measures are being introduced because it has been observed that the pictorial quality perceived by human visual system is also affected by the overall general impression of the viewed video stream on humans. In addition, recent studies have shown that human visual system awards higher response to more salient image locations and features [6].

2.2 Distortion-Invisibility

The distortion-invisibility feature measures the average amount of distortion that may be visible at each pixel with respect to a visibility threshold. The distortion-invisibility measure, $D(t)$, for each colour component of every frame of the video is given by:

$$D(t) = \left\{ 1 / \frac{1}{WH} \sum_{x=1}^W \sum_{y=1}^H \left[\gamma_1 + \frac{\hat{d}(x, y, t)}{\gamma_2 + T(x, y, t)} \right] \right\}$$

$T(x, y, t)$ is the visibility threshold at a particular pixel location (x, y) and time interval t , W and H are width and height of the

video frame respectively, γ_1 and γ_2 are included to prevent possible division by zero in the equation. Also,

$$\hat{d}(x, y, t) = \begin{cases} 0 & \text{if } d(x, y, t) \leq (1-s) \cdot T(x, y, t) \\ d(x, y, t) & \text{otherwise} \end{cases}$$

where s is the soft-criterion to avoid clipping of the difference data near the visibility threshold T , and $d(x, y, t)$ is the difference between a frame in the test video I_d and the reference video I_o at the same pixel location (x, y) and time t and is defined as:

$$d(x, y, t) = |I_o(x, y, t) - I_d(x, y, t)|$$

The visibility threshold T is given by:

$$T(x, y, t) = \left(T^l(x, y, t) + T^s(x, y, t) - C^{ls} \cdot \min\{T^l(x, y, t), T^s(x, y, t)\} \right) T^m(x, y, t)$$

The visibility threshold $T(x, y, t)$ provides an indication of the maximum allowable distortions at a particular pixel in the image frame which will still not be visible to human eyes. Here, $T^l(x, y, t)$, $T^s(x, y, t)$ and $T^m(x, y, t)$ can be regarded as effects due to colour masking, spatial-textural masking, and temporal masking respectively at a particular pixel located at position (x, y) in the image frame at time interval t in the video sequence, while C^{ls} is a constant term which accounts for the overlapping effect in masking. Masking is a very important visual phenomenon which explains why similar artefacts are disturbing in certain regions of an image frame while they are hardly noticeable in other regions.

The temporal masking T^m attempts to emulate the effect of human vision's characteristic of being able to accept higher video-frame distortion due to larger temporal changes:

$$T^m(x, y, t) = e^{f_r \cdot f_s} \begin{cases} T_1^m & \text{if } |d_f(x, y, t)| \leq T_3^m \\ T_1^m \left(\frac{1 - ((L_m - d_f(x, y, t)) / (L_m - T_3^m))}{z_2} - 1 \right) + T_2^m & \text{if } d_f(x, y, t) < -T_3^m \\ T_o^m \left(\frac{1 - ((L_m - d_f(x, y, t)) / (L_m - T_3^m))}{z_1} - 1 \right) + T_2^m & \text{Otherwise} \end{cases}$$

where $d_f(x, y, t)$ is the inter-frame difference at a particular pixel location (x, y) in time t between a current frame $I_o(x, y, t)$ and a previous coded frame $I_o(x, y, t - f_f/f_r)$ (assuming that frames that has been coded at below full frame rate has been repeated in this video sequence) and is mathematically expressed as:

$$d_f(x, y, t) = I_o(x, y, t) - I_o(x, y, t - f_f/f_r)$$

Here, f_r is the frame rate at which the video has been compressed, f_s is a scaling factor, while L_m , T_o^m , T_1^m , T_2^m , T_3^m , z_1 , and z_2 are constants used to determine the exact profile of the temporal masking.

The colour masking T^l attempts to emulate the effect of human vision's characteristic of being able to accept higher video-frame distortion when the background colour is above or below a certain mid-level threshold:

$$T^l(x, y, t) = \begin{cases} T_1^l \left(\frac{1 - ((L_l/2 + b(x, y, t)) / (L_l - r))}{v_2} - 1 \right) + T_2^l & \text{if } b(x, y, t) \leq (L_l - r) \\ T_o^l \left(\frac{1 - ((3L_l/2 - b(x, y, t)) / (L_l - r))}{v_1} - 1 \right) + T_2^l & \text{if } b(x, y, t) > (L_l + r) \\ T_2^l & \text{Otherwise} \end{cases}$$

Here, T_o^l , T_1^l and T_2^l , L_l , r , v_1 , and v_2 are constants used to determine the exact profile of the colour masking.

The spatial-textural masking T^s attempts to emulate the effect of human vision's characteristic of being able to accept higher video-frame distortion when the particular point has richer texture or spatial profile:

$$T^s(x, y, t) = \left(\frac{m(x, y, t)b(x, y, t)\alpha_1 + m(x, y, t)\alpha_2}{b(x, y, t)\alpha_3 + \alpha_4} \right) W(x, y, t)$$

Here, α_1 , α_2 , α_3 , and α_4 are constants used to determine the exact profile of the spatial-textural masking.

In the spatial-textural masking, $m(x, y, t)$ is the weighted average colour $g_k(x, y)$ in four different orientations (weighted by weight values w_1 and w_2) and it attempts to capture the textural-masking characteristic of the small local region centred on pixel (x, y, t) and can be defined as:

$$m(x, y, t) = \frac{1}{4} |w_1 g_1(x, y, t) + w_2 g_2(x, y, t) + w_2 g_3(x, y, t) + w_1 g_4(x, y, t)|$$

The different weightings ($w_1 = 1.5$ and $w_2 = 0.5$) given to the horizontal and vertical directions and the diagonal directions are for taking into consideration the difference in sensitivity of the human visual system to different spatial orientations. The psycho-visual data available from past research shows that the sensitivity of the human visual system is orientation dependent. However, it has been found that the sensitivity to vertical orientations is similar to that of horizontal ones [19]. In an experiment which examined the sensitivity of the HVS to frequencies at various spatial orientations [20], it was found that HVS could better resolve frequencies along the horizontal and vertical orientation than along the diagonal orientations.

Also, $g_k(x, y, t)$ is the average colour around a pixel located at position (x, y) of a frame in the original reference video sequence at time interval t and is computed by convolving a 7x7 mask, G_k , with this particular frame in the original reference video sequence:

$$g_k(x, y, t) = \sum_{m=-3}^3 \sum_{n=-3}^3 f(x+m, y+n, t) \cdot G_k(m+4, n+4, t)$$

The four 7x7 masks, G_k , for $k=\{1, 2, 3, 4\}$, are four differently oriented gradient masks used to capture the strength of the gradients around a pixel located at position (x, y, t) .

Here, $b(x, y, t)$ is the average background colour around a pixel located at position (x, y) of a frame in the original reference video sequence at time interval t and is computed by convolving a 7x7 low-pass filter mask, B , with this particular frame in the original reference video sequence:

$$b(x, y, t) = \sum_{m=-3}^3 \sum_{n=-3}^3 f(x+m, y+n, t) \cdot B(m+4, n+4, t)$$

In addition, $W(x, y, t)$ is an edge-adaptive weight of the pixel at location (x, y) of a frame in the original reference video sequence at time interval t , and it attempts to reduce the spatial-textural masking at edge locations because artefacts that are found on essential edge locations tend to reduce the visual quality of the image frame. Previous research findings have reported that edge information is found to be of primary importance in visual perception [11, 8]. Edge is directly related to the image content that demarcates object boundaries, surface crease, and other important visual events. Distortion at an edge is easier to be

noticed than that in other textured regions because edge structure attracts more visual attention from human visual system [8]. Thus, for more accurate visibility threshold estimation, spatial-textural masking in edge and non-edge regions has to be distinguished, as adopted here.

The corresponding edge-adaptive weight matrix W , obtained by convolving \hat{E} with a 7x7 low-pass filter g , is given by:

$$W = \hat{E} * g; \quad \hat{E} = 1 - (0.9E)$$

where $*$ is a convolution operator, E is the edge matrix of the original image frame computed with any edge detection technique and contains values of 1 and 0 for edge and non-edge pixels respectively.

2.3 Block-Fidelity

The block-fidelity feature measures the amount of distortion at block-boundaries in the test video when compared to the original reference (undistorted) video. The blocking effect is one of the significant coding artefacts that often occur in video compression. The block-fidelity measure for each colour component of each individual frame of the video is defined as:

$$F_{BF}(t) = e^{(0.25) \left(\left| \left(B_d^h(t) + B_d^v(t) \right) - \left(B_o^h(t) + B_o^v(t) \right) \right| \right) / \left(B_o^h(t) + B_o^v(t) \right)}$$

where the subscript o refers to the original video sequence, d refers to the test video sequence, and:

$$B^h(t) = \frac{1}{H(\lfloor W/4 \rfloor - 1)} \sum_{y=1}^H \sum_{x=1}^{\lfloor W/4 \rfloor - 1} |d^h(4x, y, t)|$$

$$d^h(x, y, t) = I(x+1, y, t) - I(x, y, t)$$

$I(x, y, t)$ denotes the colour value of the input image frame I at pixel location (x, y) and time interval t , H is the height of the image, W is the width of the image, $x \in [1, W]$, $y \in [1, H]$, $t \in [1, N]$, and N is the total number of frames in the video sequence. $B^v(t)$ and $d^v(x, y, t)$ can be computed in a similar way except in the y -direction.

2.4 Content Richness Fidelity

The content richness fidelity feature measures the fidelity of the richness of test video's content when compared to the reference video. This feature closely correlates with human perceptual response which tends to assign better subjective ratings to more lively and colourful images.

The image content richness fidelity feature for each colour component of every individual frame of time interval t of the video can be defined as:

$$F_{RF}(t) = e^{(0.25)R_d(t)/R_o(t)}; \quad R(t) = - \sum_{p(i) \neq 0} p(i) \log_e(p(i))$$

$$p(i) = N(i) / \sum_{\forall i} N(i)$$

Here, i is the colour value, $N(i)$ is the number of occurrence of i in the image frame, and $p(i)$ is the probability or relative frequency of i .

3 PERFORMANCE OF PROPOSED METRIC

3.1 Test Conditions

Ninety test video sequences are generated by subjecting 12 different original undistorted CIF and QCIF video sequences ("Container", "Coast Guard", "Japan League", "Foreman", "News", and "Tempete") to H.264 video compression with different bit rates (from 24 kbps to 384 kbps) and frame rates (from 7.5Hz to 30Hz). The bit rates under test are much lower than those used in [14] after the image size factor has been offset. Each of the video sequence consists of 250 frames.

3.2 Subjective Test Method

The subjective video quality tests of the test video sequences have been carried out as the tests conducted for the evaluation of JVT video sequences [1], using Double-Stimulus Impairment Scale variant II (DSIS-II) subjective test method and performed by 20 subjects. The decoded sequences with frame rate lower than 30 fps are displayed with repeated frames on the 30 Hz display device.

3.3 Performance

Performance is measured by comparing the metric output Q with the subjective rating of subjective tests between the original and the distorted sequences. To facilitate monotonicity of prediction and a common analysis space of comparison, Q is fitted via a 4-parameter cubic polynomial (as in [14, 13]) to the corresponding subjective rating as:

$$q = a_0 + a_1(Q) + a_2(Q^2) + a_3(Q^3)$$

Two performance measures have been used for comparison here (as in [14, 13]): (1) Pearson correlation coefficient (r_p), and (2) Spearman rank-order correlation coefficient (r_s). In the ideal match between a metric's outputs and subjective ratings, $r_p = 1$ and $r_s = 1$.

Table 1 shows the results of the proposed metric with respect to PSNR (The parameters of the proposed metric have been obtained empirically based on a small subset of the data set). The upper bound (UB) and lower bound (LB) of Pearson correlation were obtained with a confidence interval of 95%. It can be seen that the proposed video quality metric (with a Pearson correlation of 0.916 and Spearman correlation of 0.920) performs much better than the PSNR (which has a Pearson correlation of 0.701 and Spearman correlation of 0.676).

Figure 1 shows the scatterplot of subjective ratings versus the PSNR values, while Figure 2 shows the scatterplot of subjective ratings versus the video quality ratings estimated using our proposed metric. In these two figures, the middle solid line portrays the logistic fit using the above-mentioned 4-parameter cubic polynomial, while the upper dotted curve and the lower dotted curve portray the upper bound and lower bound respectively obtained with a confidence interval of 95%.

| | r_p | r_p UB | r_p LB | r_s |
|-----------------|-------|----------|----------|-------|
| PSNR | 0.701 | 0.793 | 0.578 | 0.676 |
| Proposed metric | 0.916 | 0.944 | 0.875 | 0.920 |

Table 1: Performance of proposed video quality metric and PSNR

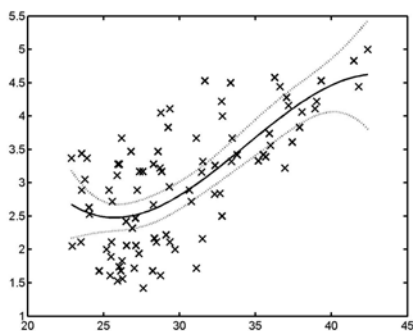


Figure 1: Scatterplot of subjective ratings vs PSNRs

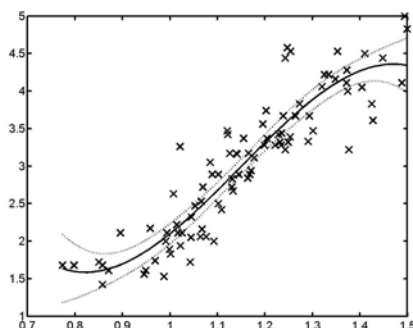


Figure 2: Scatterplot of subjective ratings vs outputs of proposed metric

5. CONCLUSION

In this paper, an objective video quality metric to automatically assess the perceived quality of a stream of video images has been described. The proposed method attempts to emulate human visual perception by introducing computational models based on block-fidelity, content richness fidelity, spatial-textural, colour, and temporal maskings. This method has been tested on digitally coded CIF and QCIF video sequences (at 24~384 Kbps and 30~7.5Hz) and shown to achieve significantly better correlation with subjective viewing results in comparison with the PSNR measure. Such an objective video quality metric will be extremely useful as it can replace the use of performance measure such as the traditionally used PSNR which has been found to correlate poorly with subjective quality ratings and also subjective tests which is not only time-consuming but also tedious and expensive to perform.

6. REFERENCES

- [1] Baroncini, V., ISO/IEC JTC 1/SC29/WG 11, 4240, Sidney, July 2001.
- [2] Chou, C.H., and Li, Y.C., "A perceptually tuned subband image coder based on the measure of just-noticeable-distortion profile", *IEEE Trans. on Circuits & Systems for Video Technology*, Vol. 5(6), 1995, pp. 467-476.
- [3] Girod, B., "What's wrong with mean-squared error", *Digital Images and Human Vision*, MIT Press, 1993, pp. 207-220.
- [4] Jayant, N., Johnston, J., and Safranek, R., "Signal compression based on models of human perception", *Proceedings of IEEE*, Vol. 81, 1993, pp. 1385-1422.
- [5] Karunasekera, S.A., and Kingsbury, N.G. "A distortion measure for blocking artifacts in image based on human visual sensitivity", *IEEE Trans. on Image Processing*, Vol. 4(6), 1995, pp.713-724.
- [6] Li, Z., "A saliency map in primary visual cortex", *Trends in Cognitive Science*, Vol. 6(1), 2002.
- [7] Lindh, P., and Lambrecht, C., "Efficient spatio-temporal decomposition for perceptual processing of video sequences", *International Conference on Image Processing*, 1996.
- [8] Marr, D., *Vision: A computational investigation into the human representation and processing of visual information*, W.H. Freeman and Co., 1982.
- [9] Miyahara, M., Kotani, K. and Algazi, V.R., "Objective picture quality scale (PQS) for image coding", *IEEE Transactions on Communications*, Vol. 46(9), 1998, pp. 1215-125.
- [10] Nothdurft, H.C., "Saliency from feature contrast: Additivity across dimensions", *Vision Res.*, Vol. 40, 2000, pp. 1183-1201.
- [11] Ran, X., and Farvardin, N., "A perceptually motivated three-component image model – Part I: Description of the model", *IEEE Trans. on Image Processing*, Vol. 4(4), 1995, pp. 401-415.
- [12] Tan, K.T., and Ghanbari, M., "A multi-metric objective picture-quality measurement model for MPEG video", *IEEE Trans. on Circuits & Syst. for Video Tech*, Vol. 10(7), 2000, pp. 1208-1213.
- [13] VQEG (Video Quality Expert Group), "Evaluation of new methods for objective testing of video quality: objective test plan", ITU-T/COM-T/COM12/C, www.vqeg.org, 1998.
- [14] VQEG (Video Quality Expert Group), "Final Report from the Video Quality Expert Group on the validation of Objective Models of Video Quality Assessment", www.vqeg.org, 2000.
- [15] Watson, A.B., Hu, J., and McGowan III, J.F., "DVQ: A digital video quality metric based on human vision", *Journal of Electronic Imaging*, Vol. 10(1), 2001, pp. 20-29.
- [16] Winkler, S., "A perceptual distortion metric for digital color video", *SPIE Proc. Human Vision and Elect. Imaging IV*, Vol. 3644, pp. 175 – 184, 1999.
- [17] Wu, H.R. and Yuen, M., "A generalize block-edge impairment metric for video coding", *IEEE Signal Processing Letters*, Vol. 4(11), 1997, pp.317-320.
- [18] Yu, Z, Wu, H.R., Winkler, S. and Chen, T., "Vision-model-based impairment metric to evaluate blocking artifacts in digital video", *Proc. IEEE*, Vol. 90(1), pp. 154-169, 2002.
- [19] Campbell, F.W., and Kulikowski, J.J., "Orientation selectivity of human visual system", *Journal of Physiology*, Vol. 187, 1966, pp. 437-445.
- [20] Campbell, F.W., Kulikowski, J.J., and Levinson, J., "The effect of orientation on the visual resolution of gratings", *Journal of Physiology*, Vol. 187, 1966, 427-436.