# PERCEPTUAL REPRESENTATIONS OF STRUCTURAL INFORMATION IN IMAGES: APPLICATION TO QUALITY ASSESSMENT OF SYNTHESIZED VIEW IN FTV SCENARIO

*Suiyi Ling, Jing Li, Patrick Le Callet*

IPI/LS2N Lab, University of Nantes, France

*Junle Wang*

Turing Lab, Tencent, China

## ABSTRACT

As the immersive multimedia techniques like Free-viewpoint TV (FTV) develop at an astonishing rate, user's demand for high-quality immersive contents increases dramatically. Unlike traditional uniform artifacts, the distortions within immersive contents could be non-uniform structure-related and thus are challenging for commonly used quality metrics. Recent studies have demonstrated that the representation of visual features can be extracted from multiple levels of the hierarchy. Inspired by the hierarchical representation mechanism in the human visual system (HVS), in this paper, we explore to adopt structural representations to quantitatively measure the impact of such structure-related distortion on perceived quality in FTV scenario. More specifically, a bio-inspired full reference image quality metric is proposed based on 1) low-level **contour** descriptor; 2) mid-level contour **category** descriptor; and 3) task-oriented **non-natural structure** descriptor. The experimental results show that the proposed model outperforms significantly the state-of-the-art metrics.

***Index Terms***— Perceptual representation, structural information, image quality assessment, immersive multimedia, Free viewpoint TV

## 1. INTRODUCTION

With the rise of 3D displays, head-mounted displays and other advanced display techniques, immersive media applications such as FTV, 3DTV, Virtual Reality (VR) and LightField (LF) have become a hot topic for media ecosystems. The development of immersive media largely relies on the usage of computer vision/image processing techniques to generate synthetic contents that are likely subject to affect structures of images/videos and the viewing experience, a typical example is the synthesized virtual views in FTV scenario due to the limited camera setting/bandwidth. Quality control of the entire immersive system is thus vital for delivering acceptable quality service to users. So far, the structure-related distortions are challenging for commonly used quality metrics to quantify as they distribute locally and non-uniformly throughout the image/video. One of the best instinctive ways to predict the impacts of the non-uniform structure-related

distortions on perceived quality is to employ the representation mechanism within HVS [1].

The process of human analyzing a visual scene has been characterized by the presence of regions in the extrastriate cortex that are selectively responsive to scenes [2, 3]. These regions have often been interpreted to represent high-level properties of scenes and they also exhibit substantial sensitivity to low and mid-level properties. A recent bio-vision study [4] proposes a hierarchical framework of visual perception, which comprises a series of discrete stages that successively produce three levels of representations. This framework is illustrated in the left part of Figure 1.
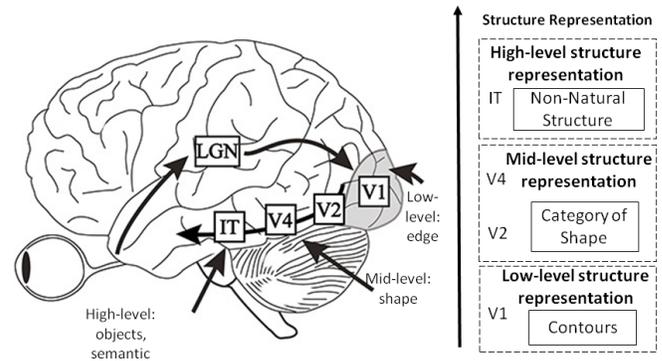


**Fig. 1**: Left: The hierarchical feedforward framework of visual perception (figure adapted from [5]). The visual precept is formed based on successive extractions and representations of low-, mid- and high-level features processed through LGN, V1, V2, V4 and IT [3, 2, 5]. Right: The proposed hierarchical model for structure-related distortion representation.

Inspired by the aforementioned bio-vision theories, in this paper, a hierarchical structure representation model is proposed and applied to quality assessment for synthesized images in FTV scenario. This model consists of three level representations as depicted in the right part of Figure 1, where low-level structure representation of images is defined as local basic structure information (e.g., local contours); mid-level is defined as intermediate 'pattern-based encoded feature', where the patterns are learned by summarizing semantic characteristics of local structure information (e.g., categories of contours); high-level is defined as 'task-related ab-

straction', which learns a set of meaningful abstract structure-related patterns reflecting the characteristics of the task (e.g., non-natural structure in image quality assessment).

The paper has the following organization: Section 2 summarizes the existing quality metrics designed for multi-view images. The proposed hierarchical metric is introduced in Section 3. The performance of the proposed metric is reported and analyzed in Section 4. Conclusions are presented in Section 5.

## 2. RELATED WORK

In order to better evaluate the quality of synthesized views in the case of FTV, some metrics are proposed. The very first metric VSQA [6] was proposed using three visibility maps which characterize complexity in terms of textures, diversity of gradient orientations and presence of high contrast. The 3DswIM was introduced by Battisti *et al.* [7] based on statistical features of wavelet sub-bands. Stanković [8] *et al.* first deployed morphological wavelet decomposition for quality assessment of synthesized images named MW-PSNR. Later, another metric devises PSNR with morphological pyramids decomposition (MP-PSNR) was proposed in [9]. Based on the fact that PSNR is more consistent with human judgment when calculated in higher morphological decomposition scales, they further proposed the reduced versions of the two metrics [10], i.e., MW-PSNR$_r$ and MP-PSNR$_r$, which provide better performance compared to the full versions. Targeting the problem that global shifting artifacts are generally over-penalized by point-wise metrics, CT-IQM [11] was proposed using an encoding scheme based on the context tree. To quantify the deformation of curves in synthesized views, an elastic metric based EM-IQM is proposed in [12]. Li *et al.* [13] proposed LOGs by considering both the geometric distortions as well as the sharpness of the images. Apart from the full reference metrics, several no reference metrics are also proposed by the community. In [14], NIQSV was proposed by hypothesizing that high-quality images are consist of flat areas separated by edges. It is then extended to NIQSV+ [15] by considering the existence of the dis-occluded regions. Recently, a novel no reference quality metric for synthesized images namely APT was proposed in [16], where the auto-regression (AR) based local image description is employed. In addition to the metrics mentioned above, we believe that there is still room to improve the performance from a perspective of bio-visual structure representation. Details of our proposal are shown in the following section.

## 3. THE PROPOSED METRIC

In this section, we propose a full-reference image quality metric based on hierarchical structure representation. The proposed framework consists of (1) a pre-processing step for structural information extraction, (2) a hierarchical feature extraction for low, mid and high-level perceptual information extraction, and (3) a pooling step for overall quality score prediction. The overall framework is shown in Figure 2.
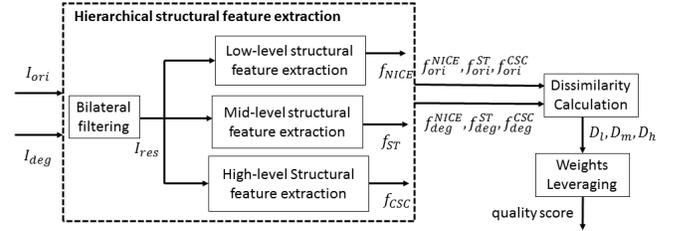


**Fig. 2**: Framework of the proposed hierarchical structural representation based image quality metric for synthesized views.

### 3.1. Structural information extraction

At the very first step of our model, we propose to separate the structural information from textural information. Previous studies [1, 17] have demonstrated that structural information plays a significantly major role in perceived quality of synthesized image compared to textural information. In addition, it has been shown that bilateral filter has the capability to emphasize such structural information [1, 17]. In our model, we adopt the approximated bilateral filter proposed in [18] for computation efficiency. Afterwards, the responses of bilateral filter for both the degraded image ($I_{deg}$) and its original one ($I_{ori}$), i.e., $f_{ori}^{NICE}$ and $f_{deg}^{NICE}$ are used as the input of the hierarchical feature extraction step. Details are shown in the following sections.

### 3.2. Low-level structure representation based estimator

As pointed out in [19] fragments of contours are the fundamental low-level structure elements that facilitate the successful identification of semantics in images. In [1, 17], it has been confirmed that the NICE (contour-based image evaluation [20]) descriptor $f^{NICE}$ plays the greatest roles in quantifying the impact of structural distortions on perceived quality, thus, it is adopted in this paper as the low-level structure representation based estimator, which is defined as

$$
\begin{aligned}
D_l &= XOR(f_{ori}^{NICE}, f_{deg}^{NICE}) \\
&= \frac{\sum^{N_c} XOR(C_{ori} \otimes E_{se}, C_{deg} \otimes E_{se})}{N_c},
\end{aligned}
\tag{1}
$$

where $C_{ori}$ and $C_{deg}$ are the contour map detected from the original and degraded images using Canny edge detector, $XOR(\cdot)$ is the point-wise exclusive-or (XOR) operation, and $N_c$ is the number of contour elements. The contour maps are subjected to morphological dilation operation (denoted as $\otimes$) with a $3 \times 3$ 'plus-sign' shaped structuring element $E_{se}$ so that the shapes within the images are probed.

### 3.3. Mid-level structure representation based estimator

HVS is very efficient in encoding the properties of stimulus by utilizing available regularities. Those efficient representations would be maximally informative with respect to the actual inputs in the world. In particular, low-level elements that share similar characteristics should be encoded more compactly [21]. A higher semantic and efficient representation of low-level structural elements, i.e., mid-level representation is thus defined, which are the categories of the contours. Based on this assumption, a Sketch-Token based Image Quality Metric (ST-IQM) by checking how the categories of contours change due to structural distortions [22] is employed as the mid-level structure representation based estimator in our study (also employed in [23] as mid-level descriptor), where contours are first 'encoded' as a vector $f^{ST}$ of contour categories likelihood values. The mid-level estimator is calculated as the Minkowski summation of the errors computed based on the mid-level descriptor across the entire image:

$$D_m = \frac{[\sum^{N_p} D_{JSD}(f_{ori}^{ST}, f_{deg}^{ST})^\beta]^{\frac{1}{\beta}}}{N_p} \qquad (2)$$

where $f_{ori}^{ST}$ and $f_{deg}^{ST}$ are the sketch-token descriptors of pairs of matched pixels from the original image to the degraded one (pixels are first matched using a registration methodology proposed in [24] to avoid over-penalizing acceptable global shifting artifacts). $D_{JSD}(\cdot)$ denotes the Jensen–Shannon divergence function, $N_p$ is the number of pixels contained in the image, and $\beta$ is a parameter corresponds to the $\beta - norm$ defining the $L^\beta$ vector space. In our study, $\beta = 4$.

### 3.4. High-level structure representation based estimator

It is mentioned in [25] that neural code in the higher-level cortex can be sparse, where each element stands for meaningful characteristics of the world (sparsity is considered as one of the essential principles to sensory representation). In [26], the process of image quality assessment is also assumed to adhere to such a strategy. For quality assessment tasks, the 'abstract' elements within the sparse dictionary could be items that reflect quality. For instance, in the case where structure-related distortions are the dominate artifacts, the items could be a set of non-natural structures. In our work, we employed a Convolutional Sparse Coding (CSC) based representation in [27, 28] as a high-level structure representation based estimator. Details are described below.

First, with a set of patches $Y$ that contain obvious local structure-related distortions collected from synthesized views, a convolutional dictionary $D_Y$ is first learned with a fast CSC algorithm proposed in [29] with the equation below

$$\underset{D_Y}{\arg\min} \frac{1}{2}\|y - \sum_{k=1}^{K} D_k \circledast Z_k\|^2 \qquad (3)$$
$$s.t. \ \|D_k\|_2^2 \leqslant 1,$$

where $\circledast$ denotes the convolution operation, $y \in Y$ denotes training samples, $Z_k$ represents sparse feature maps, $D_k$ is the $k_{th}$ convolution kernel and $K$ is the number of kernels within the dictionary.

With the learned dictionary, for a given $M \times N$ test image $I$, its sparse representation $Z_I$ could be generated with the trained dictionary:

$$\underset{Z_I}{\arg\min} \frac{1}{2}\|I - D_Y \circledast Z_I\|^2 + \alpha\|Z_I\|_1 \qquad (4)$$

where $Z_I = [Z_1; \dots; Z_k; \dots; Z_K]$ is a $M \times N \times K$ tensor of feature maps for $I$, where each map $Z_k$ is the response of using kernel $D_k$. $\alpha$ is a tunable parameter that could be used to balance the model accuracy and the sparsity of feature maps. Finally, a convolutional sparse coding based high-level feature vector $f^{CSC}$ could be then extracted for any image $I$ with:

$$f^{CSC} = (A(Z_1), ..., A(Z_K)), \qquad (5)$$

where $A(\cdot)$ is defined as

$$A(Z_k) = \frac{\sum_{i=1}^{M}\sum_{j=1}^{N} \mathbf{1}(Z_k(i,j) > \varepsilon)}{M \times N}, \qquad (6)$$

$\mathbf{1}(c)$ is an indicator function that equals to 1 if the specified binary clause $c$ is true and 0 otherwise, and $\varepsilon$ is a threshold for selecting activated pixels. Function $A(\cdot)$ aggregates the number of pixels which are above the threshold $\varepsilon$ in each sparse feature map $Z_k$ corresponding to each kernel $D_k$. Intuitively, this function counts the number of pixels that are activated by the corresponding kernel. Since the kernels are trained to capture structured-related artifacts, this process could be interpreted as the computation of certain types of relative artifacts in the entire image and thus could be used to indicate perceived quality. Finally, support vector regression (denoted as $SVR(\cdot)$) is used to predict the final quality score with the CSC based feature using 1000 times cross-validation. Here, the model that yields the median performance $m_{med}$ is used to compute the high-level structural dissimilarity $D_h$:

$$D_h = |SVR(f_{ori}^{CSC}) - SVR(f_{deg}^{CSC})| \qquad (7)$$

### 3.5. Quality score prediction

The quality score $S$ is then predicted with the linear combination of the three-level structural distortions $D_l$, $D_m$ and $D_h$ after normalization so that the dissimilarity values are in a range of [0,1]:

$$S = w_l \cdot D_l + w_m \cdot D_m + w_h \cdot D_h$$
$$s.t. \ w_l + w_m + w_h = 1, \qquad (8)$$

where $w_l$, $w_m$, and $w_h$ are the weights used for fine-tuning the roles of the low, mid and high-level structural representation based estimators respectively.

## 4. EXPERIMENTAL RESULTS

The performance of the proposed model is evaluated on the IRCCyN/IVC DIBR images database [30]. Images from this database were obtained from three multi-view video plus depth sequences: 'Book Arrival', 'Lovebird1' and 'Newspaper'. Seven DIBR algorithms processed the three sequences to generate four new virtual views for each of them. The database is composed of 84 synthesized views and 12 original frames extracted from the corresponding sequences along with subjective scores. After calculating the differential Mean Opinion Score (DMOS), the following widely employed criteria are utilized to evaluate the performances of the quality metrics: Pearson Correlation Coefficient (PCC), Spearmans rank order Correlation Coefficient (SCC) and Root Mean Squared Error (RMSE). Please note that non-linear mapping between the subjective scores and objective measures [15] is conducted before calculating the PCC, SCC, and RMSE.

### 4.1. Parameters selection

In this study, The overall performance of the proposed model is reported with a configuration of $w_l = 0.05$, $w_m = 0.25$ and $w_h = 0.75$ that obtains the median performance throughout a 1000 cross validation as described in [31, 15] with a constraint that the sum of them equals to one. To further analyze the roles of the three levels structural representation in quantifying the impact of structural distortions in predicting quality, the performances of different configurations are checked. The results are shown in Figure 3. It could be observed from the figure that the performance increases with a higher $w_h$. This observation verifies the fact that, high-level structure representation is of greater capability in quantifying structure-related distortions since it is more task-oriented.
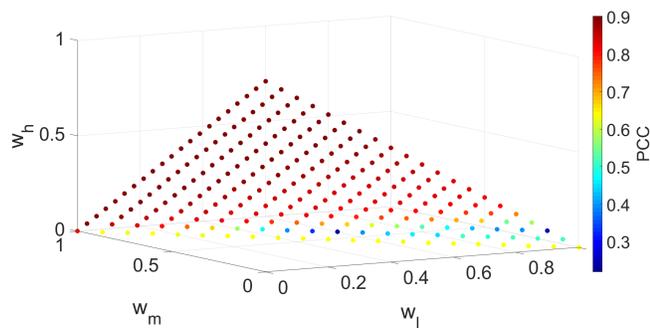


**Fig. 3**: Performances of different configurations of $w_l$, $w_m$ and $w_h$.

### 4.2. Overall performance

The overall performance results are shown in Table 1. According to the table, the proposed hierarchical structural representation based model outperforms the compared state-of-the-art FR/NR image quality metrics designed for quality as-

sessment of synthesized views in FTV scenario. It obtains gains of 9.2% and 8.6% in PCC values compared to the second best performing FR metric LoGs and the best performing NR metric CSC-NRM respectively.

**Table 1**: Performance comparison of the proposed metric with state-of-the-art metrics

|  | PCC | SCC | RMSE |
|---|---|---|---|
| Full Reference Metric (FR) | | | |
| 3DSwIM [7] | 0.6864 | 0.4842 | 0.6125 |
| MP-PSNR$_r$ [10] | 0.6954 | 0.4784 | 0.6606 |
| MW-PSNR$_r$ [10] | 0.6637 | 0.4921 | 0.6293 |
| CT-IQM [11] | 0.6809 | 0.6626 | 0.4877 |
| BF-M [1] | 0.6980 | 0.5885 | 0.4768 |
| EM-IQM [12] | 0.7430 | 0.6726 | 0.4455 |
| ST-IQM [22] | 0.8217 | 0.7710 | 0.3929 |
| LoGs [13] | 0.8256 | 0.7812 | 0.3601 |
| **Proposed** | **0.9023** | **0.8448** | **0.2870** |
| NO Reference Metric (NR) | | | |
| NIQSV [14] | 0.6346 | 0.5146 | 0.6167 |
| NIQSV+ [15] | 0.7114 | 0.4679 | 0.6668 |
| APT [16] | 0.7307 | 0.7140 | 0.4622 |
| CSC-NRM [25] | **0.8302** | **0.7827** | **0.3233** |

To analyze if the performances of the proposed metric and other well performed FR and NR metrics are significant, the F-test based on the residual difference between the predicted objective scores and the subjective DMOS values as described in [32] is employed. The result is reported in Table 2, where '1' indicates the quality metric in the row outperforms significantly the one in the column. Thus, the proposed metric outperforms the others significantly.

**Table 2**: Statistic significance results based on F-test.

| Metric | LoGs | ST-IQM | NIQSV+ | APT | CSC-NRM |
|---|---|---|---|---|---|
| Proposed | 1 | 1 | 1 | 1 | 1 |

## 5. CONCLUSION

Local, non-uniform structure-related distortions within immersive multimedia are challenging for traditional quality metrics. Inspired by the hierarchical framework of visual perception, in this paper, a 3-level structure representation based model is proposed. This model quantifies the structure-related distortion by checking 1) how local contours change (low-level); 2) how the categories of contour change (mid-level); 3) and the amount of non-natural structure within the synthetic image compared to the original image (high-level). The role of each level of representations on image quality assessment has been studied as well. According to experimental results, the proposed model is significantly superior to the state-of-the-art metrics.

# 6. REFERENCES

[1] Suiyi Ling, Patrick Le Callet, and Zitong Yu, "The role of structure and textual information in image utility and quality assessment tasks," *Electronic Imaging*, vol. 2018, no. 14, pp. 1–13, 2018.

[2] Iris IA Groen, Edward H Silson, and Chris I Baker, "Contributions of low-and high-level properties to neural processing of visual scenes in the human brain," *Phil. Trans. R. Soc. B*, vol. 372, no. 1714, pp. 20160102, 2017.

[3] Timothy J Andrews, David M Watson, Grace E Rice, and Tom Hartley, "Low-level properties of natural images predict topographic patterns of neural response in the ventral visual pathway," *Journal of Vision*, vol. 15, no. 7, pp. 3–3, 2015.

[4] Jonathan W Peirce, "Understanding mid-level representations in visual processing," *Journal of Vision*, vol. 15, no. 7, pp. 5–5, 2015.

[5] Mauro Manassi, Bilge Sayim, and Michael H Herzog, "When crowding of crowding leads to uncrowding," *Journal of Vision*, vol. 13, no. 13, pp. 10–10, 2013.

[6] Pierre-Henri Conze, Philippe Robert, and Luce Morin, "Objective view synthesis quality assessment," in *IS&T/SPIE Electronic Imaging*. International Society for Optics and Photonics, 2012, pp. 82881M–82881M.

[7] Federica Battisti, Emilie Bosc, Marco Carli, Patrick Le Callet, and Simone Perugia, "Objective image quality assessment of 3d synthesized views," *Signal Processing: Image Communication*, vol. 30, pp. 78–88, 2015.

[8] Dragana Sandić-Stanković, Dragan Kukolj, and Patrick Le Callet, "Dibr synthesized image quality assessment based on morphological wavelets," in *Quality of Multimedia Experience (QoMEX), 2015 Seventh International Workshop on*. IEEE, 2015, pp. 1–6.

[9] Dragana Sandic-Stankovic, Dragan Kukolj, and Patrick Le Callet, "Dibr synthesized image quality assessment based on morphological pyramids," in *2015 3DTV-Conference: The True Vision-Capture, Transmission and Display of 3D Video (3DTV-CON)*. IEEE, 2015, pp. 1–4.

[10] Dragana Sandić-Stanković, Dragan Kukolj, and Patrick Le Callet, "Dibr-synthesized image quality assessment based on morphological multi-scale approach," *EURASIP Journal on Image and Video Processing*, vol. 2017, no. 1, pp. 4, 2016.

[11] Patrick Le Callet Ling, Suiyi and Cheung Gene, "Quality assessment for synthesized view based on variable-length context tree," in *Multimedia Signal Processing (MMSP), 2017 IEEE 19th International Workshop on*. IEEE, 2017.

[12] Suiyi Ling and Patrick Le Callet, "Image quality assessment for dibr synthesized views using elastic metric," in *Proceedings of the 2017 ACM on Multimedia Conference*. ACM, 2017, pp. 1157–1163.

[13] Leida Li, Yu Zhou, Ke Gu, Weisi Lin, and Shiqi Wang, "Quality assessment of dibr-synthesized images by measuring local geometric distortions and global sharpness," *IEEE Transactions on Multimedia*, vol. 20, no. 4, pp. 914–926, 2018.

[14] Shishun Tian, Lu Zhang, Luce Morin, and Olivier Deforges, "Niqsv: A no reference image quality assessment metric for 3d synthesized views," in *Acoustics, Speech and Signal Processing (ICASSP), 2017 IEEE International Conference on*. IEEE, 2017, pp. 1248–1252.

[15] Shishun Tian, Lu Zhang, Luce Morin, and Olivier Déforges, "Niqsv+: A no-reference synthesized view quality assessment metric," *IEEE Transactions on Image Processing*, vol. 27, no. 4, pp. 1652–1664, 2018.

[16] Ke Gu, Vinit Jakhetiya, Jun-Fei Qiao, Xiaoli Li, Weisi Lin, and Daniel Thalmann, "Model-based referenceless quality metric of 3d synthesized images using local image description," *IEEE Transactions on Image Processing*, 2017.

[17] Yashas Rai, Ahmed Aldahdooh, Suiyi Ling, Marcus Barkowsky, and Patrick Le Callet, "Effect of content features on short-term video quality in the visual periphery," in *Multimedia Signal Processing (MMSP), 2016 IEEE 18th International Workshop on*. IEEE, 2016, pp. 1–6.

[18] Sylvain Paris and Frédo Durand, "A fast approximation of the bilateral filter using a signal processing approach," *International journal of computer vision*, vol. 81, no. 1, pp. 24–52, 2009.

[19] Jamie Shotton, Andrew Blake, and Roberto Cipolla, "Multiscale categorical object recognition using contour fragments," *IEEE transactions on pattern analysis and machine intelligence*, vol. 30, no. 7, pp. 1270–1281, 2008.

[20] David M Rouse, Romuald Pépion, Sheila S Hemami, and Patrick Le Callet, "Image utility assessment and a relationship with image quality assessment," in *Human Vision and Electronic Imaging XIV*. International Society for Optics and Photonics, 2009, vol. 7240, p. 724010.

[21] Jonas Kubilius, Johan Wagemans, and Hans P Op de Beeck, "Encoding of configural regularity in the human visual system," *Journal of Vision*, vol. 14, no. 9, pp. 11–11, 2014.

[22] Suiyi Ling and Patrick Le Callet, "Image quality assessment for free viewpoint video based on mid-level contours feature," in *Multimedia and Expo (ICME), 2017 IEEE International Conference on*. IEEE, 2017, pp. 79–84.

[23] Yu Zhou, Leida Li, Suiyi Ling, and Patrick Le Callet, "Quality assessment for view synthesis using low-level and mid-level structural representation," *Signal Processing: Image Communication*, 2019.

[24] Jignesh N Sarvaiya, Suprava Patnaik, and Salman Bombaywala, "Image registration by template matching using normalized cross-correlation," in *Advances in Computing, Control, & Telecommunication Technologies, 2009. ACT'09. International Conference on*. IEEE, 2009, pp. 819–822.

[25] Peter Foldiak, "Sparse coding in the primate cortex," *The handbook of brain theory and neural networks*, 2003.

[26] Ayyoub Ahar, Adriaan Barri, and Peter Schelkens, "From sparse coding significance to perceptual quality: A new approach for image quality assessment," *IEEE Transactions on Image Processing*, vol. 27, no. 2, pp. 879–893, 2018.

[27] Suiyi Ling, Gene Cheung, and Patrick Le Callet, "No-reference quality assessment for stitched panoramic images using convolutional sparse coding and compound feature selection," in *2018 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2018, pp. 1–6.

[28] Suiyi Ling and Patrick Le Callet, "How to learn the effect of non-uniform distortion on perceived visual quality? case study using convolutional sparse coding for quality assessment of synthesized views," in *2018 25th IEEE International Conference on Image Processing (ICIP)*. IEEE, 2018, pp. 286–290.

[29] Michal Šorel and Filip Šroubek, "Fast convolutional sparse coding using matrix inversion lemma," *Digital Signal Processing*, vol. 55, pp. 44–51, 2016.

[30] Emilie Bosc, Romuald Pepion, Patrick Le Callet, Martin Koppel, Patrick Ndjiki-Nya, Muriel Pressigout, and Luce Morin, "Towards a new quality metric for 3-d synthesized view assessment," *IEEE Journal of Selected Topics in Signal Processing*, vol. 5, no. 7, pp. 1332–1343, 2011.

[31] Suiyi Ling, Jesús Gutiérrez, Ke Gu, and Patrick Le Callet, "Prediction of the influence of navigation scan-path on perceived quality of free-viewpoint videos," *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*, 2019.

[32] Xiangkai Liu, Yun Zhang, Sudeng Hu, Sam Kwong, C-C Jay Kuo, and Qiang Peng, "Subjective and objective video quality assessment of 3D synthesized views with texture/depth compression distortion," *IEEE Transactions on Image Processing*, vol. 24, no. 12, pp. 4847–4861, Dec. 2015.