

# Learned Scanpaths Aid Blind Panoramic Video Quality Assessment

## Supplementary Material

Kanglong Fan<sup>1</sup>, Wen Wen<sup>1</sup>, Mu Li<sup>2\*</sup>, Yifan Peng<sup>3</sup>, and Kede Ma<sup>1</sup>

<sup>1</sup> City University of Hong Kong, <sup>2</sup> Harbin Institute of Technology, Shenzhen

<sup>3</sup> The University of Hong Kong

### 1. Details of the Density Estimation Network

The architecture of the density estimation network is depicted in Figure S1. Of particular interest is the masked computation [S3] used to process the causal path. We refer the readers to the code implementation at <https://github.com/kalofan/AutoScanpathQA> for detailed parameter configurations.

### 2. Relative $uv$ Coordinate System

Figure S2 compares different coordinate systems. The transformation  $(u, v) = \Psi_t(\phi, \theta)$  that maps the viewpoint  $(\phi, \theta)$  to the  $(u, v)$  coordinates relative to the  $t$ -th viewport centered at  $(\phi_t, \theta_t)$  is broken down into multiple steps. First,  $(\phi, \theta)$  is transformed to  $(x, y, z)$  in the Cartesian coordinate system:

$$x = r \cos(\phi) \cos(\theta), \quad (\text{S1})$$

$$y = r \cos(\phi) \sin(\theta), \quad (\text{S2})$$

$$z = r \sin(\phi), \quad (\text{S3})$$

where  $r = 0.5W_v \cot(0.5\theta_v)$  is the radius of the sphere, determined by the width of the viewport,  $W_v$  and the field of view,  $\theta_v$ . Second,  $(x, y, z)$  is rotated with respect to  $(\phi_t, \theta_t)$ :

$$\begin{pmatrix} x_t \\ y_t \\ z_t \end{pmatrix} = R(\phi_t, \theta_t)^\top \times \begin{pmatrix} x \\ y \\ z \end{pmatrix}. \quad (\text{S4})$$

where  $R(\phi_t, \theta_t)$  is the rotation matrix defined as the product of two matrices  $R_2 \times R_1$ :

$$R_1 = \begin{pmatrix} a & -b & 0 \\ b & a & 0 \\ 0 & 0 & 1 \end{pmatrix}, \quad (\text{S5})$$

$$R_2 = \begin{pmatrix} c + (1-c)b^2 & -(1-c)ab & -da \\ -(1-c)ab & c + (1-c)a^2 & -db \\ da & db & c \end{pmatrix}, \quad (\text{S6})$$

\*Corresponding author.

where

$$a = \cos(\theta_t), \quad (\text{S7})$$

$$b = \sin(\theta_t), \quad (\text{S8})$$

$$c = \cos(\phi_t), \quad (\text{S9})$$

$$d = \sin(\phi_t). \quad (\text{S10})$$

After rotation,  $(x_t, y_t, z_t)$  is transformed to  $(r, y'_t, z'_t)$  by projecting it to the plane  $x = r$ . Last,  $(r, y'_t, z'_t)$  is readily represented in the  $uv$  coordinate system:

$$u = y'_t + 0.5W_v - 0.5, \quad (\text{S11})$$

$$v = 0.5H_v - z'_t - 0.5, \quad (\text{S12})$$

where  $H_v$  is the height of the viewport. We may shift the origin to the viewport center, leading to

$$u = y'_t, \quad (\text{S13})$$

$$v = -z'_t. \quad (\text{S14})$$

### 3. PID Controller

The PID controller [S1] is a prevalent feedback mechanism that enables continuous modulation of control signals. We adopt the sampling strategy of [S3] and assume a proxy viewer governed by Newton's laws of motion. Initially, the proxy viewer is positioned at the starting point  $\hat{\mathbf{r}}_{-1} = (0, 0)$  in the  $uv$  coordinate system, with an initial speed  $\mathbf{b}_{-1}$  and acceleration  $\mathbf{a}_{-1}$ . We determine the  $t$ -th viewpoint by

$$\hat{\mathbf{r}}_t = \hat{\mathbf{r}}_{t-1} + \Delta\lambda\mathbf{b}_{t-1} + \frac{1}{2}(\Delta\lambda)^2\mathbf{a}_{t-1}, t \in \{0, \dots, W-1\}, \quad (\text{S15})$$

where the speed  $\mathbf{b}_{t-1}$  is updated by

$$\mathbf{b}_t = \mathbf{b}_{t-1} + \Delta\lambda\mathbf{a}_{t-1}. \quad (\text{S16})$$

$\Delta\lambda$  represents the sampling interval, *i.e.*, the inverse of the sampling rate. To update the acceleration  $\mathbf{a}_{t-1}$ , a reference viewpoint  $\tilde{\mathbf{r}}_t$  for  $\hat{\mathbf{r}}_t$  is provided by sampling from

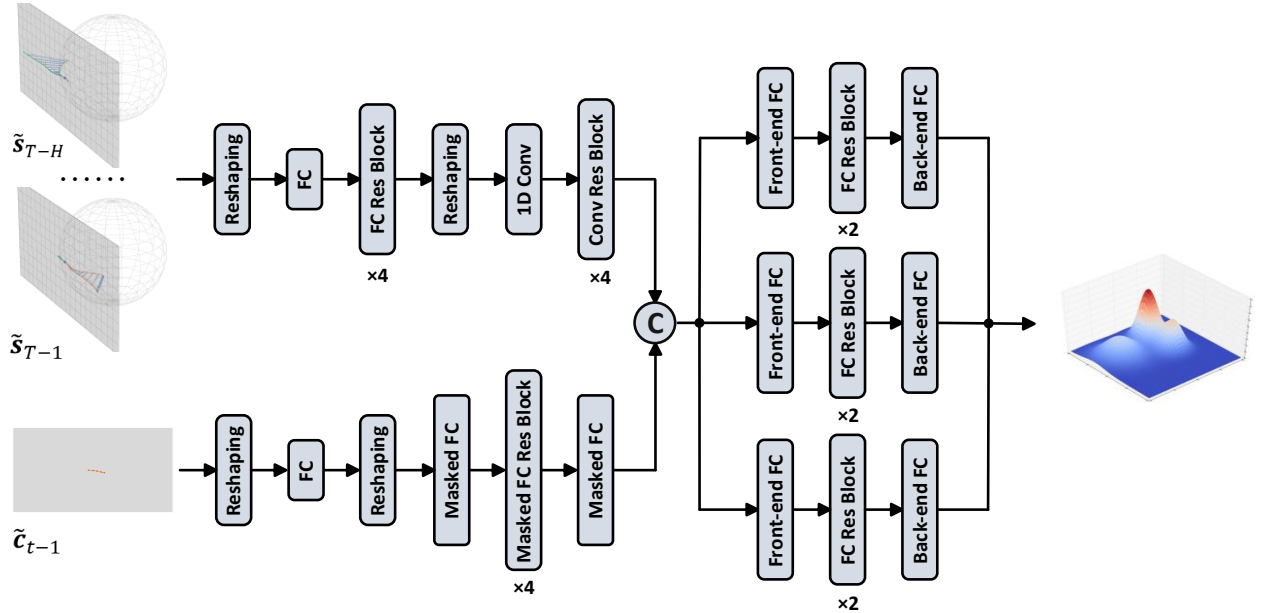


Figure S1. Architecture of the density estimation network.

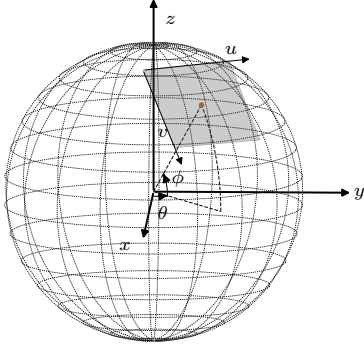


Figure S2. Illustration of different coordinate systems relevant to panoramic signal processing: Spherical Euler coordinates  $(\phi, \theta)$ , 3D Euclidean coordinates  $(x, y, z)$ , and  $uv$  coordinates  $(u, v)$ . Image adapted from [S3].

$P(\tilde{r}_t | s, c_t)$ , where  $c_t = \{\hat{r}_0, \dots, \hat{r}_{t-1}\}$ . Consequently, an error signal is generated:

$$e_t = \tilde{r}_t - \hat{r}_t, \quad (\text{S17})$$

which is fed to the PID controller for acceleration adjustment:

$$a_t = K_p e_t + K_i \sum_{\tau=0}^t e_\tau + K_d (e_t - e_{t-1}), \quad (\text{S18})$$

where  $K_p$ ,  $K_i$ , and  $K_d$  denote the proportional, integral, and derivative gains, respectively. During training, we

back-propagate the gradient through the PID controller to the scanpath generator for parameter update.

#### 4. Implementation Details

The first stage of training is carried out by the Adam method [S2] on VRVQW with an initial learning rate of  $10^{-4}$  and a minibatch size of 48. After the 50-th epoch, the learning rate decays by a ratio of 0.1, and we pre-train the scanpath generator for a total of 100 epochs. The parameters for the PID controller are determined using the Ziegler–Nichols method [S5].

For the second and third stages of training, the Adam method is also employed, and the detailed settings of different quality assessors are as follows.

**ScanpathVQA.** In the second training stage, the quality assessor is trained for 30 epochs, with an initial learning rate of  $5 \times 10^{-5}$ , a decay ratio of 0.95 per 2 epochs, and a batch size of 8. In the third stage of training on VRVQW, the entire method is trained for 5 epochs, with an initial learning rate of  $10^{-6}$ , a decay ratio of 0.1 after the 2-nd epoch, and a batch size of 4. In the third stage of training on CVIQD and OIQA, the method is trained for 5 epochs, with an initial learning rate of  $10^{-5}$ , a decay ratio of 0.9 per epoch, and a batch size of 4.

**GSR-S/GSR-X.** In the second training stage, we follow the settings described in the original paper [S4]. In the third stage, the entire method is trained for 10 epochs, with an initial learning rate of  $10^{-6}$ , a decay ratio of 0.9 per 2 epochs, and a batch size of 8. The input configuration also follows

the original paper [S4].

**Assessor360.** The settings in the second training stage are identical to those of ScanpathVQA. In the third stage, the entire method is trained for 10 epochs, with an initial learning rate of  $10^{-5}$ , a decay ratio of 0.9 per 2 epochs, and a batch size of 4. Owing to limitations in computer memory, the number of viewport sequences  $N$  is reduced from 20 to 15.

## References

- [S1] Kiam Heong Ang, Gregory Chong, and Yun Li. PID control system analysis, design, and technology. *IEEE Transactions on Control Systems Technology*, 13(4):559–576, 2005. [1](#)
- [S2] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015. [2](#)
- [S3] Mu Li, Kanglong Fan, and Kede Ma. Scanpath prediction in panoramic videos via expected code length minimization. *arXiv preprint arXiv:2305.02536*, 2023. [1](#), [2](#)
- [S4] Xiangjie Sui, Hanwei Zhu, Xuelin Liu, Yuming Fang, Shiqi Wang, and Zhou Wang. Perceptual quality assessment of 360° images based on generative scanpath representation. *arXiv preprint arXiv:2309.03472*, 2023. [2](#), [3](#)
- [S5] John G. Ziegler and Nathaniel B. Nichols. Optimum settings for automatic controllers. *Transactions of the American Society of Mechanical Engineers*, 64(8):759–765, 1942. [2](#)