

# High-fidelity Face Tracking for AR/VR via Deep Lighting Adaptation

Lele Chen<sup>1,2</sup>   Chen Cao<sup>1</sup>   Fernando De la Torre<sup>1</sup>   Jason Saragih<sup>1</sup>   Chenliang Xu<sup>2</sup>   Yaser Sheikh<sup>1</sup>  
<sup>1</sup> Facebook Reality Labs   <sup>2</sup> Univeristy of Rochester

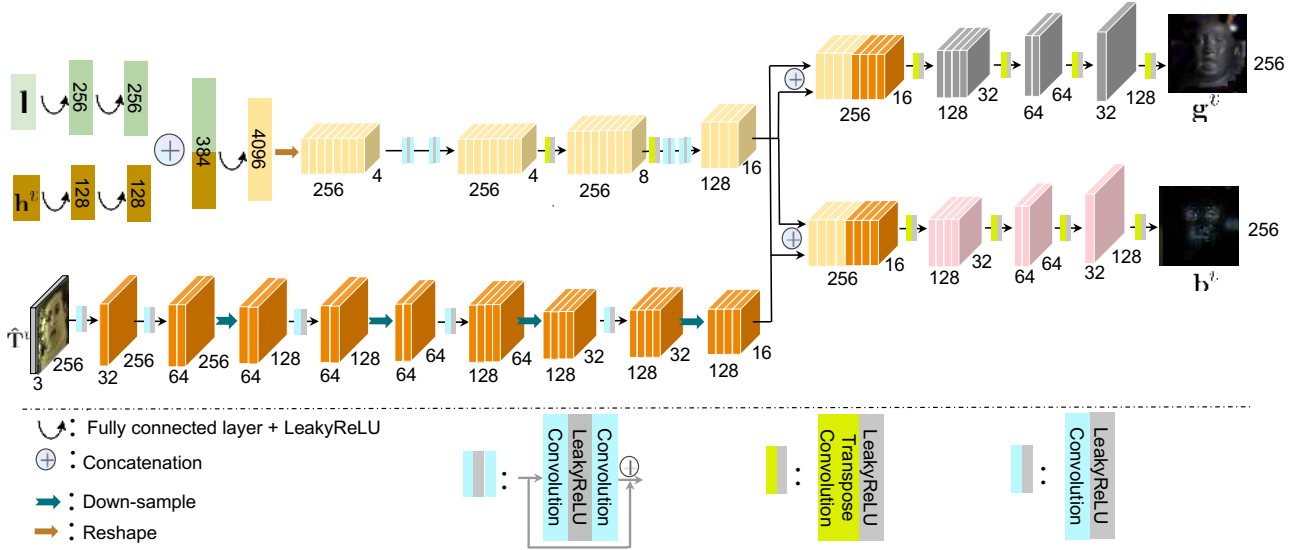


Figure 1. The detailed network structure of our lighting transfer network ( $G$ ).

In this supplementary file, we explain the network details of our lighting model.

1024 by bilinear interpolation.

## A. Network Structure

We present the detailed network structure in Fig. 1.

## B. Inputs and Outputs

The lighting code  $\mathbf{l}$  is a pre-defined vector when we train  $G$  on light-stage data, and is a learnable vector when we refine  $G$  on in-the-wild video frames. During training on light-stage data, the lighting direction is encoded by the position of the non-zero element in  $\mathbf{l}$ , and the lighting color is encoded by the value of the non-zero element in  $\mathbf{l}$ . The view-dependent head pose  $\mathbf{h}^v \in \mathbb{R}^6 = \{\mathbf{r}_x, \mathbf{r}_y, \mathbf{r}_z, \mathbf{v}_x^v, \mathbf{v}_y^v, \mathbf{v}_z^v\}$ , where  $\mathbf{r} = \{\mathbf{r}_x, \mathbf{r}_y, \mathbf{r}_z\}$  and  $\mathbf{v}^v = \{\mathbf{v}_x^v, \mathbf{v}_y^v, \mathbf{v}_z^v\}$  are rigid head rotation and viewpoint vector, respectively. The fully-lit texture  $\hat{\mathbf{T}}^v$  is obtained from DAM decoder, and we down-sample it to the size of  $3 \times 256 \times 256$ .

The outputs are the gain and bias map  $\mathbf{g}^v, \mathbf{b}^v$ , and we upsample the output  $\mathbf{g}^v, \mathbf{b}^v$  back to the size of  $3 \times 1024 \times$