# Total-Editing: Head Avatar with Editable Appearance, Motion, and Lighting

## Supplementary Material

### A. Training Scheme and Objective Functions

We begin by pre-training our lightmap estimator using synthetic data with ground truths. The objective function is

$$\mathcal{L}_{\text{pre}} = \mathcal{L}_{\mathbf{S}} + \mathcal{L}_{\mathbf{N}},\tag{22}$$

where  $\mathcal{L}_{\mathbf{S}}$  is an L1 loss comparing the predicted and ground truth lightmaps  $\mathbf{S}_d$ ,  $\{\mathbf{S}_s(n)\}$ ,  $\mathcal{L}_{\mathbf{N}}$  is a cosine similarity loss between predicted and ground truth normal  $\mathbf{N}$ . After pretraining, we detach the normal decoder  $\mathcal{E}_{nor}$  and integrate the rest of the lightmap estimator with Total-Editing. The entire network is then trained end-to-end using both real and synthetic data. During training, we randomly sample appearance source  $\mathbf{I}_{app}$  and motion source  $\mathbf{I}_{mot}$  of one subject, with the editing target  $\mathbf{I}^*$  equal to the motion source  $\mathbf{I}_{mot}$ . As for the lighting source, we use the HDR environment map of  $\mathbf{I}_{mot}$  for synthetic data and another random frame from the same video clip for real data. In this phase, our reconstruction objective is

$$\mathcal{L}_{rec} = \mathcal{L}_1 + \mathcal{L}_{LPIPS} + \mathcal{L}_{id} + \mathcal{L}_{seg} + \mathcal{L}_{\mathbf{a}} + \mathcal{L}_{\mathbf{n}} + \mathcal{L}_{\mathbf{S}}, \ (23)$$

where  $\mathcal{L}_1$ ,  $\mathcal{L}_{LPIPS}$ ,  $\mathcal{L}_{id}$  are pixel-wise L1, perceptual difference [80], and negative cosine similarity of face recognition features [14] between the editing result  $\hat{\mathbf{I}}$  and the target  $\mathbf{I}^*$ ,  $\mathcal{L}_{seg}$  and  $\mathcal{L}_{a}$  are L1 losses for the rendered foreground mask and albedo,  $\mathcal{L}_{n}$  is a cosine similarity loss for the rendered normal, and  $\mathcal{L}_{S}$  is the L1 loss for estimated lightmaps. Note that  $\mathcal{L}_{a}$  and  $\mathcal{L}_{S}$  are used only for synthetic data, while  $\mathcal{L}_{n}$  is also applied to real data with Sapiens [33] pseudo ground truths. Further, we introduce regularization

$$\mathcal{L}_{\text{reg}} = \mathcal{R}_{\text{TV}} + \mathcal{R}_{\delta} + \mathcal{R}_{\mathbf{n}}, \tag{24}$$

where  $\mathcal{R}_{TV}$  is the total variation loss to promote spatial smoothness,  $\mathcal{R}_{\delta}$  is a L1 regularization for residual color  $\delta \mathbf{c}$  which constraints it from dominating the render, and

$$\mathcal{R}_{\mathbf{n}} = \left\| 1 - \mathbf{n} \cdot \left( -\frac{\nabla \sigma(\mathbf{p})}{\|\nabla \sigma(\mathbf{p})\|_2} \right) \right\|_1 \tag{25}$$

regularizes normal **n** to align with the unit negative gradient of density  $\sigma$ . We also apply an adversarial loss  $\mathcal{L}_{adv}$  with a dual discriminator [8]. Finally, the training objective is

$$\mathcal{L} = \mathcal{L}_{rec} + \mathcal{L}_{reg} + \mathcal{L}_{adv}. \tag{26}$$

### **B.** More Qualitative Results

We present additional face reenactment results on the VFHQ [71] and HDTF [81] datasets in Figs. 11 to 14,



Figure 10. Samples from one identity of the Lumos [75] dataset. We refer to each column as a unique subject, since they have different appearances and accessories.

demonstrating the effectiveness of Total-Editing in both motion and lighting control. Further, we explore two downstream applications. As shown in Fig. 15, Total-Editing enables relighting of animated portraits using HDR environment maps, producing a background replacement effect. In Fig. 16, Total-Editing can leverage arbitrary portrait images as lighting sources, vividly transferring the illumination effect from one portrait to another.

#### C. Dataset Details

Our synthetic dataset is designed to advance research in general portrait editing, offering two primary subsets: a multi-view subset and a video-like subset.

- The multi-view subset comprises 50K subjects, each captured in two distinct environments and viewed from 10 camera angles. This subset provides extensive data for analyzing objects from diverse perspectives and ensuring multi-view consistency. Samples are shown in Fig. 17.
- The video-like subset includes 10K subjects, each rendered across 10 environments with varied poses and expressions, making it well-suited for studying motion and temporal changes. Samples are shown in Fig. 18.

With diverse samples demonstrated in Fig. 19, our synthetic dataset consists of 50K subjects and 2M images. It is enriched with ground truth albedo, normal, depth, UV maps, segmentation masks, and HDR environment maps. This dataset addresses critical limitations compared to the existing synthetic datasets. For example, as shown in Fig. 10, the Lumos dataset [75] captures each subject only from one view, limiting it to tasks like portrait relighting. In contrast, our dataset incorporates multiple viewpoints and subject movements, better simulating real-world spatiotemporal variations. These improvements make our dataset more versatile and effective for downstream applications requir-

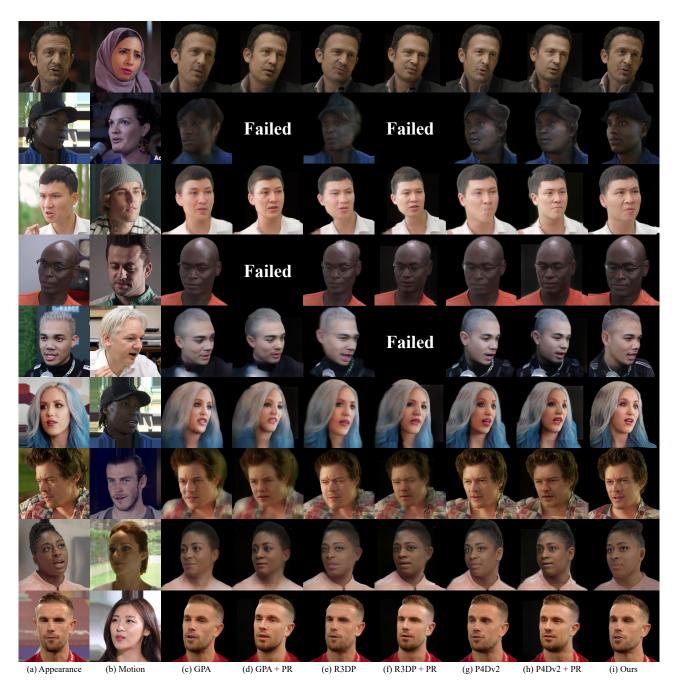


Figure 11. Additional cross-reenactment results on the VFHQ dataset.

ing spatial and/or temporal coherence.

## **D.** Evaluation Details

In Tab. 1, we exclude both generated and ground truth backgrounds from metric calculations, focusing solely on the quality of the portrait regions. Since the relighting method proposed by Cai et al. [6] requires additional cropping for

input images, we recompose the outputs with original inputs in Fig. 7 for visual consistency. During quantitative comparisons in Tab. 1, we only consider the valid areas after cropping. In Fig. 9, the backgrounds of 2<sup>nd</sup> to 3<sup>rd</sup> columns are synthesized by inpainting the lighting source portrait using [62], while those in the 4<sup>th</sup> to 5<sup>th</sup> columns are rendered from the corresponding HDR environment maps.

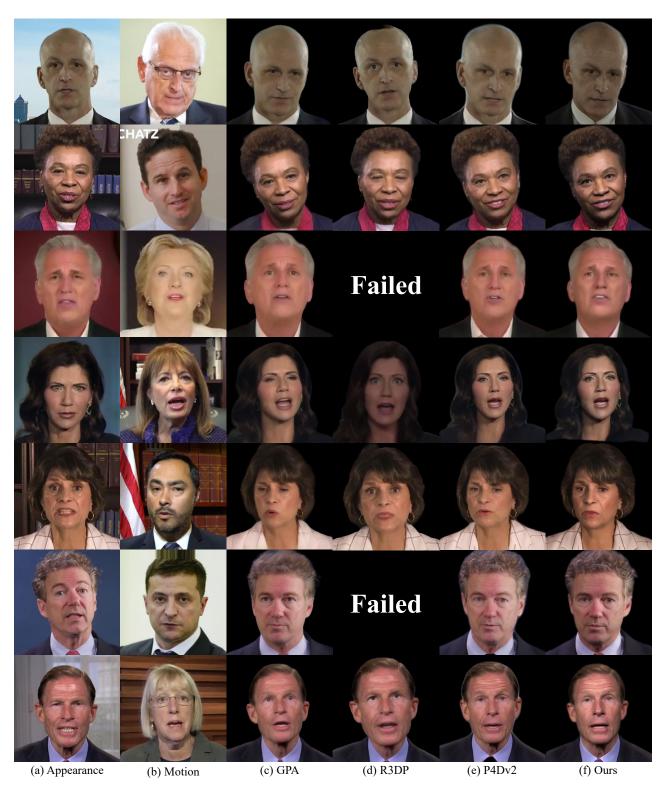


Figure 12. Cross-reenactment results on the HTDF dataset.



 $Figure\ 13.\ \textbf{Cross-reenactment\ results\ on\ the\ HTDF\ dataset\ (continued).}$ 

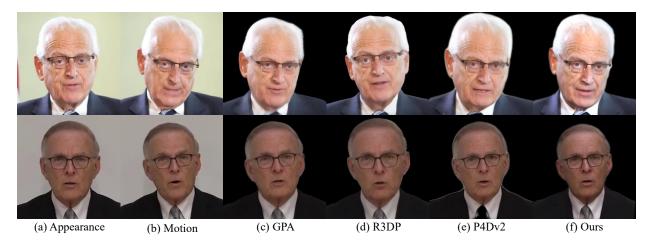


Figure 14. Self-reenactment results on the HTDF dataset.

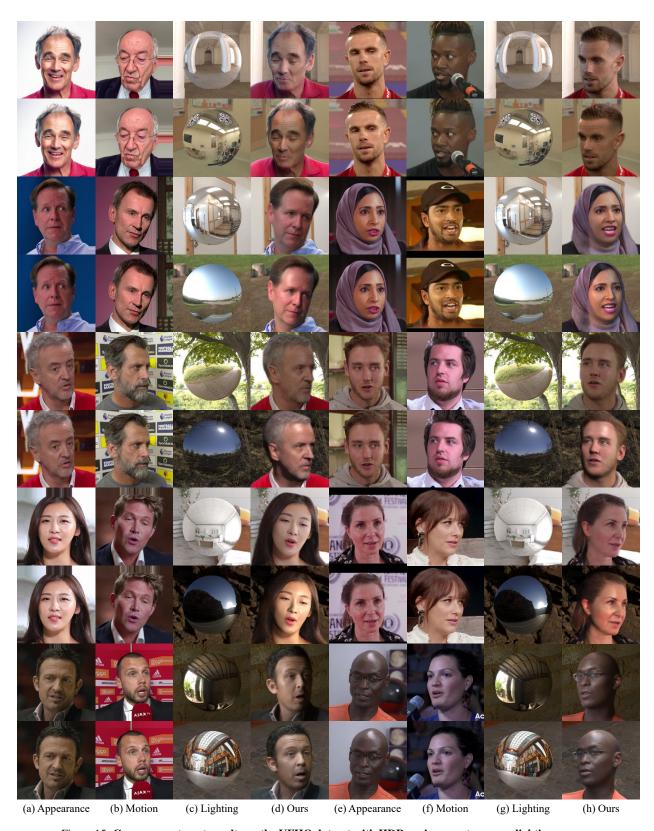


Figure 15. Cross-reenactment results on the VFHQ dataset with HDR environment maps as lighting sources.

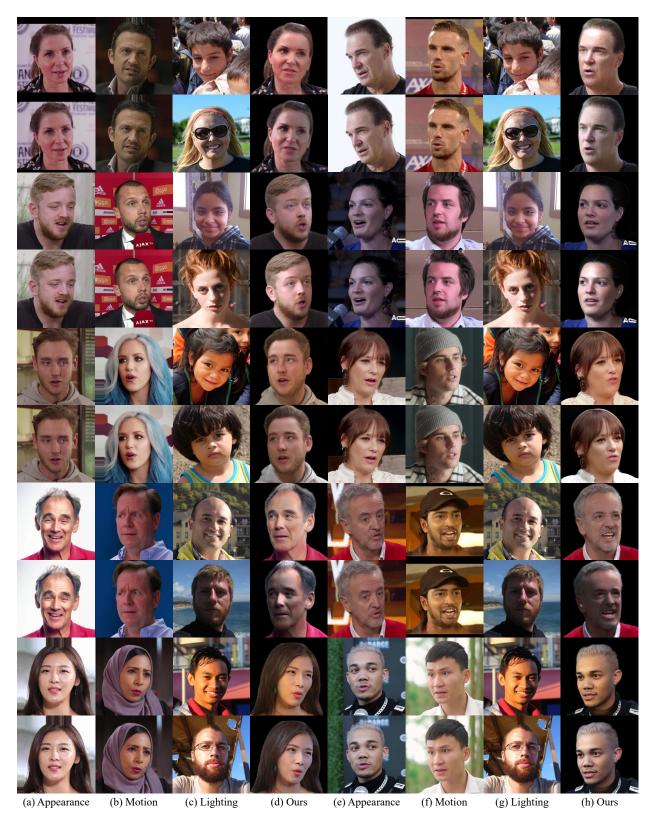


Figure 16. Cross-reenactment results on the VFHQ dataset with portrait images as lighting sources.

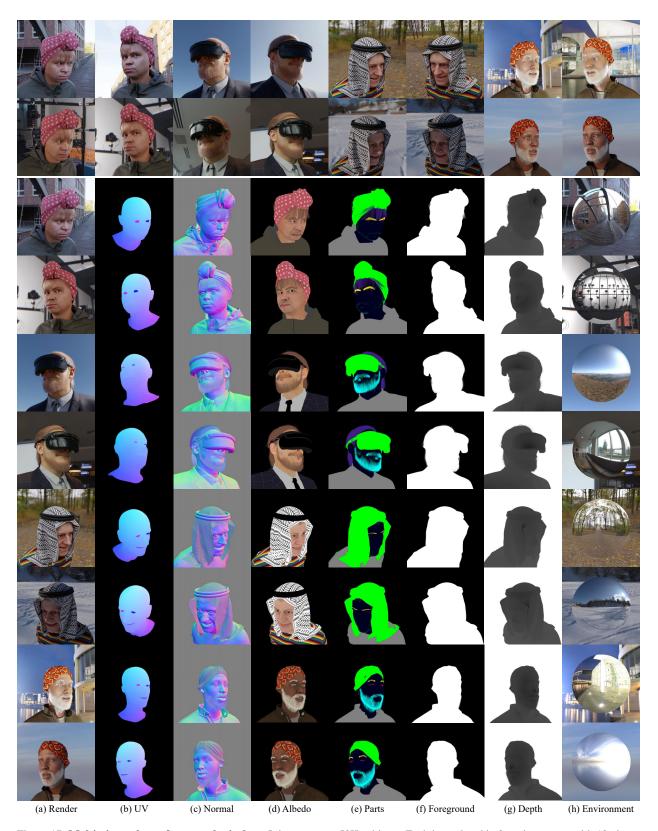


Figure 17. Multi-view subset of our synthetic data. It incorporates 50K subjects. Each is rendered in 2 environments with 10 views.

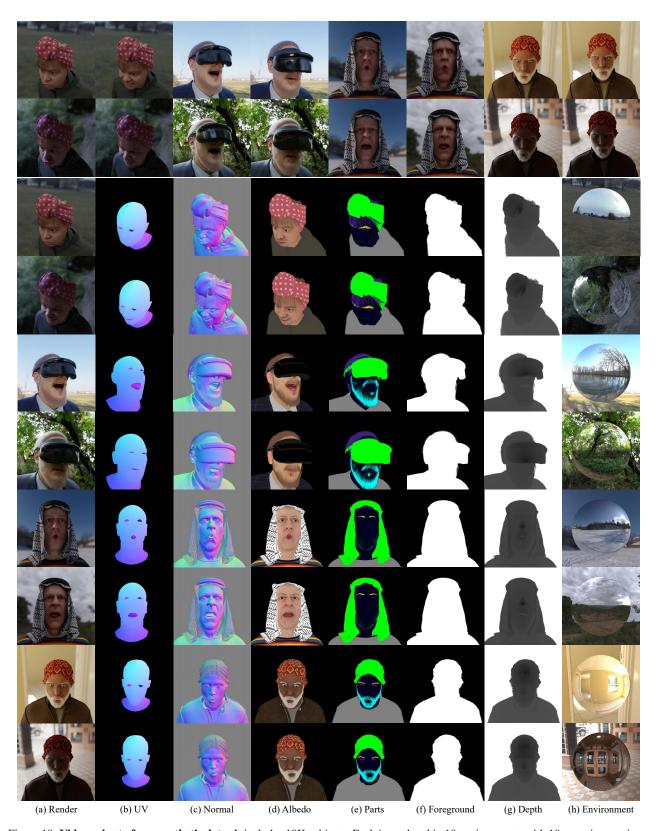


Figure 18. Video subset of our synthetic data. It includes 10K subjects. Each is rendered in 10 environments with 10 poses/expressions.

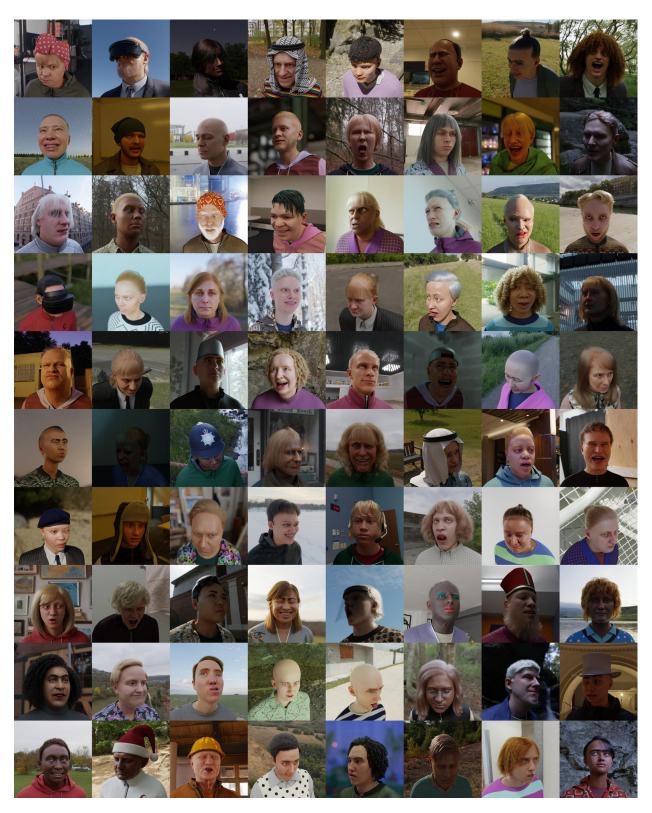


Figure 19. More subjects in our synthetic data. Subjects are with randomized poses, expressions, hairstyles, skin types, accessories, etc.