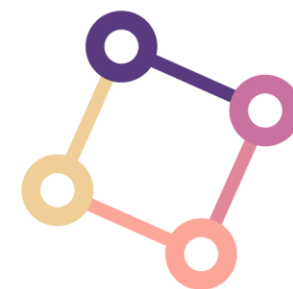# ROME: REALISTIC ONE-SHOT MESH-BASED HEAD AVATARS

## Presenter: Jaeseong Lee

ECCV Under-Review

(Released to Arxiv 17/June)

Samsung AI Moscow

DAVIAN

Data and Visual Analytics Lab

# Contents

- Preliminaries

- Overview

- Results

- Limitations

# Preliminaries – Taxonomy

Tencent AI Lab

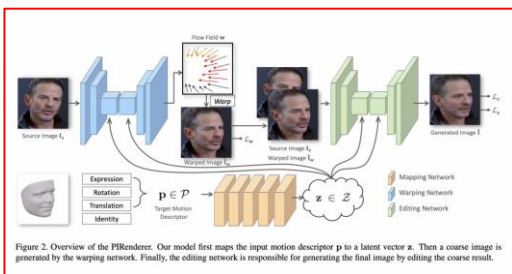Head Reenactment
(a.k.a. Neural Talking Head)

w/ 3DMM ← 2D-oper. → 3D-oper.

PIRenderer(ICCV 2021)

KP-based

No Geometry
(latent-based)

KP-based

Mesh-based

Few-shot TH(ICCV 2019)

Samsung Research

Fast-bi(ECCV 2020)

Samsung Research

LPD(CVPR 2020)

NVIDIA

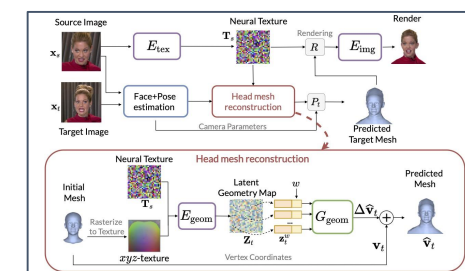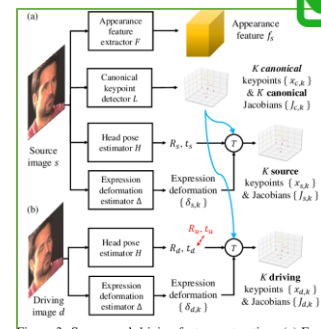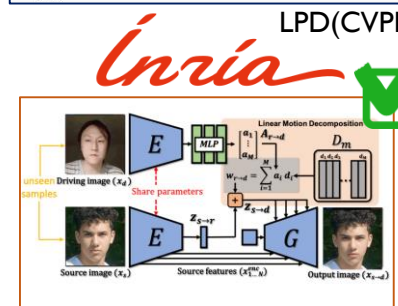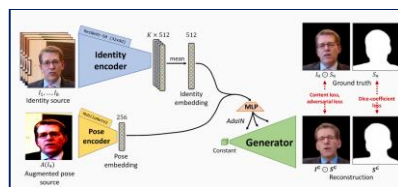Samsung Research

Snap Inc.

FOMM(NIPS 2019)

Ínría

LIA(ICLR 2022)

OSFV(CVPR 2021)

ROME(ECCV 2022(?))

# OVERVIEW(CONT'D) - STRENGTHS

- Strengths from two Perspectives

  - Head recon. : Viable to handle non-facial parts(e.g., hear and torso)



* From DECA(SIGGRAPH 2021)

  - Talking head : Viable to handle unseen facial parts or large drv/src-discrepancy (e.g., one-side facing smthng)



| Source image | Source depth | Driving image | Driving depth | FOMM [26] | OSFV [31] |

** From DAGAN(CVPR 2022)

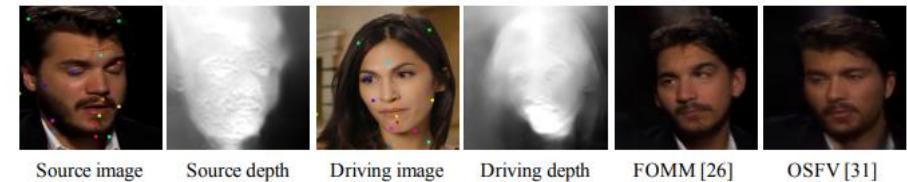# Overview – Pipeline(cont'd)

1. Rule-based data process



From off-the-shelf 3DMM extractor(DECA)

* 3DMM: representing human head from 3D scanned mesh PCA-based parameterized(Shape and Expression)

Coarse Mesh

Fine Mesh

Implies which 2d coordinate corresponds to vertex

2. How 3DMM works? How to reconstruct the Initial Mesh?



$$\boldsymbol{v}(\phi, \psi, \theta) = \text{W}(\boldsymbol{v}_{base} + \beta\phi + D\psi, \theta)$$

Where,
$\phi$: Shape parameter
$\psi$: Expression parameter
$\theta$: Head pose parameter

*Src: Shape
*Tgt: Expression and Cam/Head pose

6

# Overview – Loss functions



**Novel losses**

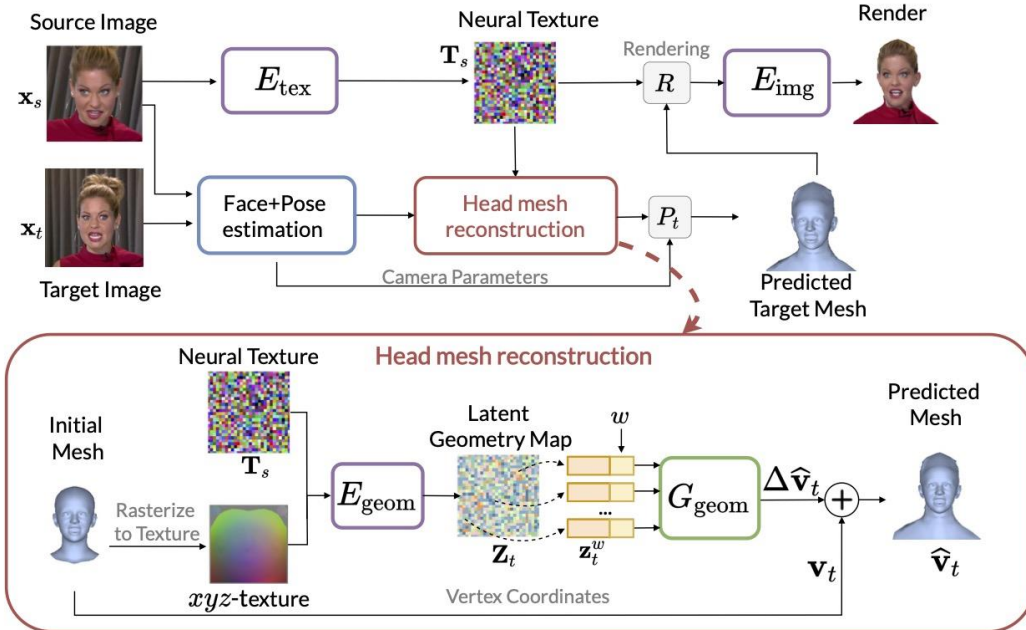$$\mathcal{L}_{\text{occ}} = \lambda_{\text{hair}} \left\| \hat{\mathbf{o}}_t^{\text{hair}} - \mathbf{s}_t^{\text{hair}} \right\|_2^2 + \lambda_{\text{o}} \left\| \hat{\mathbf{o}}_t - \mathbf{s}_t \right\|_2^2.$$

$$\mathcal{L}_{\text{chm}} = \frac{1}{2N_t} \sum_{\hat{p}_t \in \hat{\mathbf{p}}_t} \left\| \hat{p}_t - \arg\min_{p \in \mathbf{p}_t} \left\| p - \hat{p}_t \right\| \right\| +$$

$$\frac{1}{2N_t} \sum_{p_t \in \mathbf{p}_t} \left\| p_t - \arg\min_{\hat{p} \in \hat{\mathbf{p}}_t} \left\| \hat{p} - p_t \right\| \right\|.$$

$$\mathcal{L}_{\text{lap}} = \frac{1}{V} \sum_{i=1}^{V} \left\| \Delta\hat{\mathbf{v}}_i - \frac{1}{\mathcal{N}(i)} \sum_{j \in \mathcal{N}(i)} \Delta\hat{\mathbf{v}}_j \right\|_1,$$

align

regularization

Where,
$p_t$ : sampled set of 2D points in the predicted segmentation mask $\boldsymbol{s}_t$
$\hat{p_t}$ : projected 2D vertices coordinate at the target image

**Widely-used losses**

$$L_{adv} + L_{perc} + L_{arcface} + L_{dice}$$

Comparison with Head Reenactment models



| | Source | Driver | FOMM | Bi-Layer | FLAMETex | ROME |

| Method | | self-reenactment | | | cross-reenactment | |
|---|---|---|---|---|---|---|
| | LPIPS↓ | SSIM↑ | PSNR↑ | FID↓ | CSIM↑ | IQA↑ |
| FOMM | 0.09 | 0.87 | 25.8 | 52.95 | 0.53 | 55.9 |
| Bi-Layer | 0.08 | 0.83 | 23.7 | 51.4 | 0.56 | 50.48 |
| TPSMM | 0.09 | 0.85 | 26.1 | 49.27 | 0.57 | 59.5 |
| ROME | 0.08 | 0.86 | 26.2 | 45.32 | 0.62 | 66.3 |

Ablation Studies



| Input | Full | w/o $\Delta\hat{v}$ | w/o $\vec{n}$ | w/o $\mathcal{L}_{occ}$ | w/o $\mathcal{L}_{lap}$ | w/o $\mathcal{L}_{chm}$ |

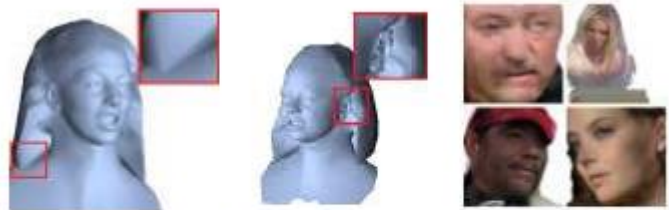Source    Driver    ROME    Source    Driver    ROME

# Limitations

- Mesh resolution



Long hair   Ear cover   Failed renders

(b) Limitations