

# StyleRig : Rigging StyleGAN for 3D Control over Portrait Images

Tewari, Ayush, et al.

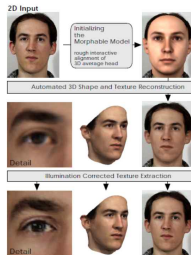
Presenter: Sejik Park

# Contents

- Introduction
- Model
- Results

# Introduction

- Improve StyleGAN controllability
  - control with 3 dimensional morphable face models (3DMMs)
    - ex. expressions, pose, illumination
  - 3DMMs
    - using 3D laser scan data
    - multivariate normal distribution
    - subregions (eyes, nose, mouth, surrounding)
    - markedness regression : neutral - expression
    - illumination-corrected texture extraction



# Model

- RigNet
  - Improve controllability by latent code modification
  - Model  
StyleGAN + differentiable face reconstruction (DFR) + RigNet
  - Train  
Pretrained StyleGAN + pretrained DFR + RigNet  
latent reconstruction loss + two-way cycle consistency losses
  - Data  
combining up 5 latent codes by StyleGAN

# Model

- RigNet

$\mathbf{I}_w$  Image correspond to latent code  $w$

$\mathbf{p}_w$  Semantic control parameters  
(input of differentiable renderer)

$\mathbf{L}_{I_w}$  Landmarks

## RigNET loss

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{rec}} + \mathcal{L}_{\text{edit}} + \mathcal{L}_{\text{consist}} .$$



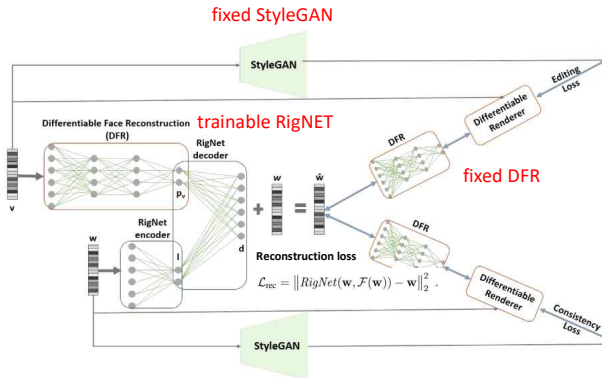
$$\mathcal{L}_{\text{edit}} = \mathcal{L}_{\text{render}}(\mathbf{I}_v, \mathbf{p}_{\text{edit}}) .$$

$$\mathcal{L}_{\text{render}}(\mathbf{I}_w, \mathbf{p}) = \mathcal{L}_{\text{photo}}(\mathbf{I}_w, \mathbf{p}) + \lambda_{\text{land}} \mathcal{L}_{\text{land}}(\mathbf{I}_w, \mathbf{p}) .$$

$$\mathcal{L}_{\text{photo}}(\mathbf{I}_w, \mathbf{p}) = \|\mathbf{M} \odot (\mathbf{I}_w - \mathcal{R}(\mathbf{p}))\|_2^2 .$$

$$\mathcal{L}_{\text{land}}(\mathbf{I}_w, \mathbf{p}) = \|\mathbf{L}_{I_w} - \mathbf{L}_M\|_2^2 ,$$

$$\mathcal{L}_{\text{consist}} = \mathcal{L}_{\text{render}}(\mathbf{I}_w, \mathbf{p}_{\text{consist}}) .$$



# Results

- Differentiable Face Reconstruction



Figure 3: Differentiable Face Reconstruction. Visualized are (image, reconstruction) pairs. The network however, only gets the latent vector corresponding to the images as input.

# Results

- Distribution of latent at different resolution and face param
  - Data : Flickr\_HQ dataset
    1. less data around Z-axis
    2. neutral or smiling/laughing faces

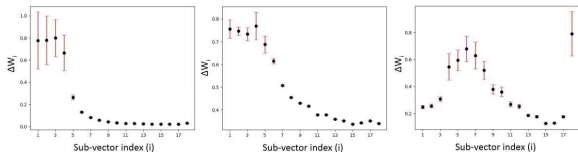


Figure 4: Change of latent vectors at different resolutions. Coarse vectors are responsible for rotation (left), medium for expressions (middle), medium and fine for illumination (right).

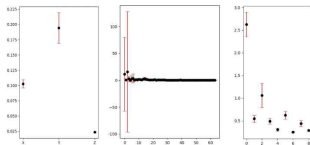


Figure 6: Distribution of face model parameters in the training data. x-axis shows the face model parameters for rotation, expression and illumination from left-right. y-axis shows the mean and variance of the parameters computed over 20k training samples.

# Results

- Mixing source & target

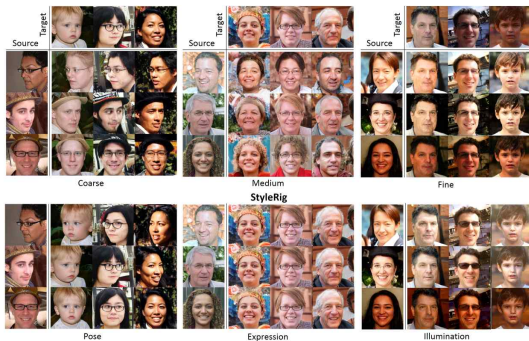


Figure 5: Mixing between source and target images generated by StyleGAN. For StyleGAN, the latent vectors of the source samples (rows) are copied to the target vectors (columns). StyleRig allows us to mix semantically meaningful parameters, i.e., head pose, expressions and scene illumination. These parameters can be copied over from the source to target images.

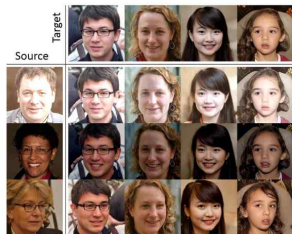
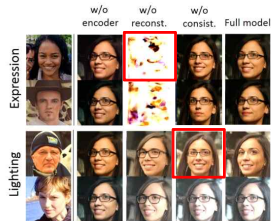
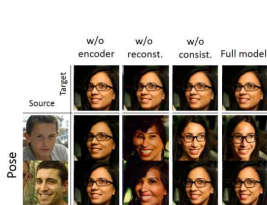
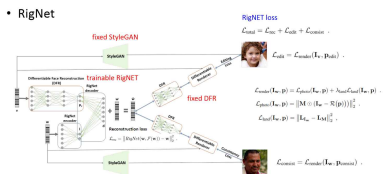


Figure 9: RigNet can also control pose, expression, and illumination parameters simultaneously. These parameters are transferred from source to target images, while the identity in the target images is preserved.



## Results

- Baseline Comparison
  - Different Loss Function
    - off reconstruction loss : non-face images
    - off consistency : one change with other feature



# Discussion

- Data
  - not able to exploit the full expressivity of the parametric face model
- Model
  - quality limitance : differentiable face reconstruction  
ex. no fine-scale detail
  - preserve parts : background, hair style
  - cf. train time (based Nvidia Volta GPU)
    - 24 hours to train StyleRig
    - more than 41 days to train StyleGAN

- End -