



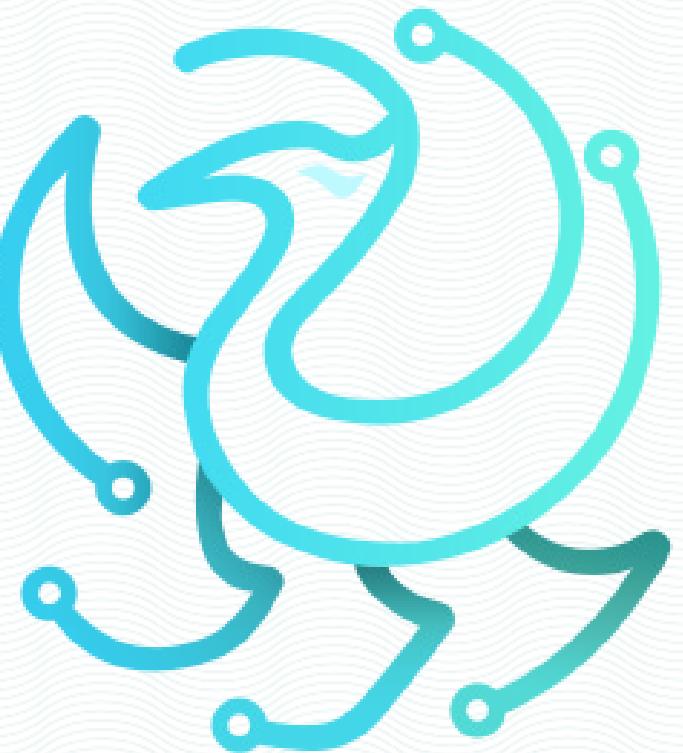
Face-to-Parameter Translation via Neural Network Renderer

Tianyang Shi

Senior AI Researcher @ Netease Games

Renjie Li

Head of AI @ Netease Games



FUXI-Lab

Established in Sep. 2017

First Game AI Research Lab in China

Vision: Enlighten Games with AI



AI

Reinforcement
Learning

User Profiling

Computer
Vision &
Graphics

Virtual
Human

Natural
Language
Processing

Big Data & Computing Platform

Talk

The AI knows you better than yourself: mining player pursuits in Justice Online
– AI Summit

Full body animation generation for expressive NPCs
– Machine Learning Summit

Building an intelligent game testing system in Netease MMORPG games
– Machine Learning Summit

Face-to-parameter translation via neural network renderer
- Core

Applying reinforcement learning to develop game AI in Netease games
- Core

Overview

- Introduction
- Face-to-Parameter Translation
- Discussion





Introduction



GAME DEVELOPERS CONFERENCE
MARCH 16–20, 2020 | #GDC20

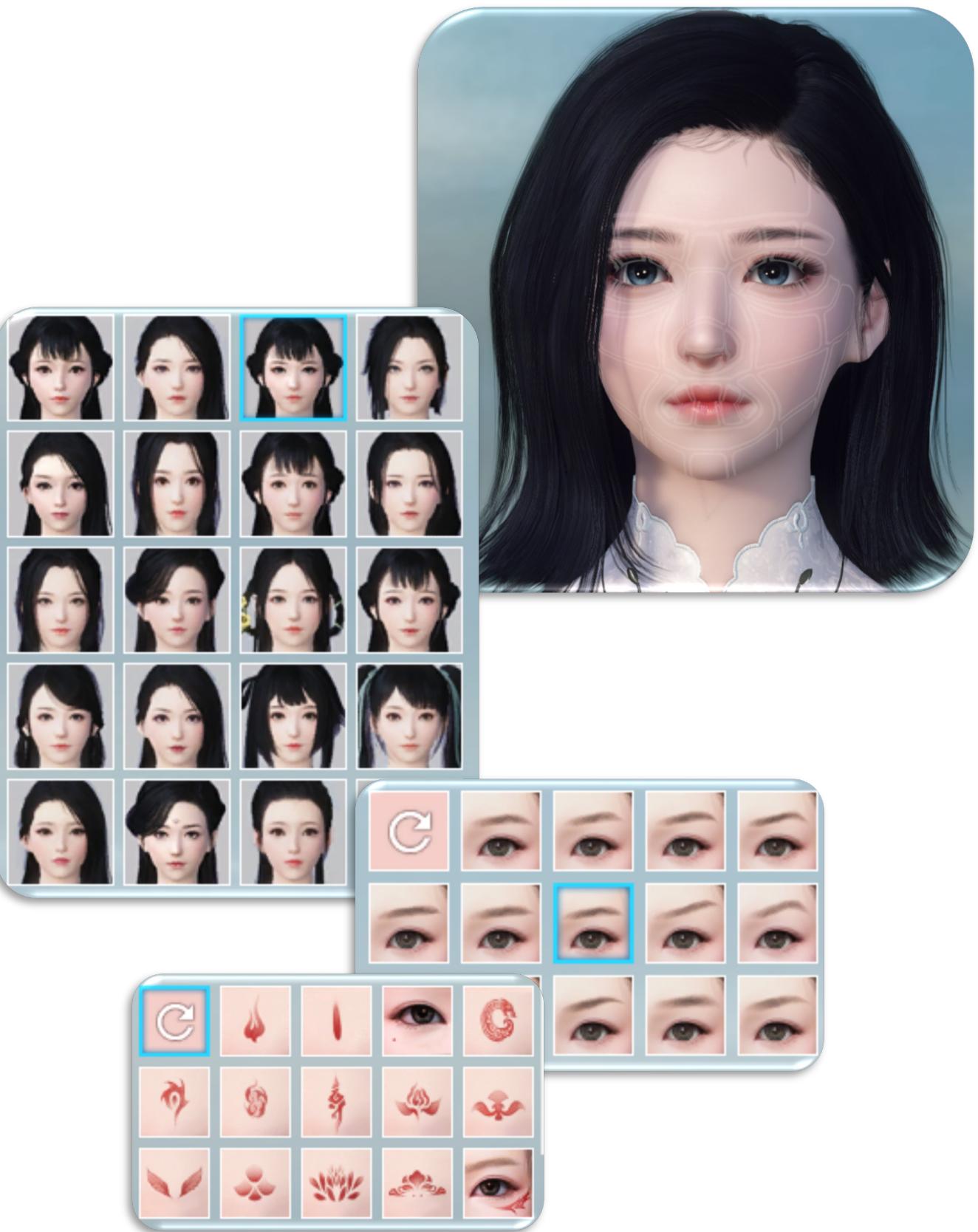
Game Character Customization

- This system allows players to edit the profiles of their in-game characters.



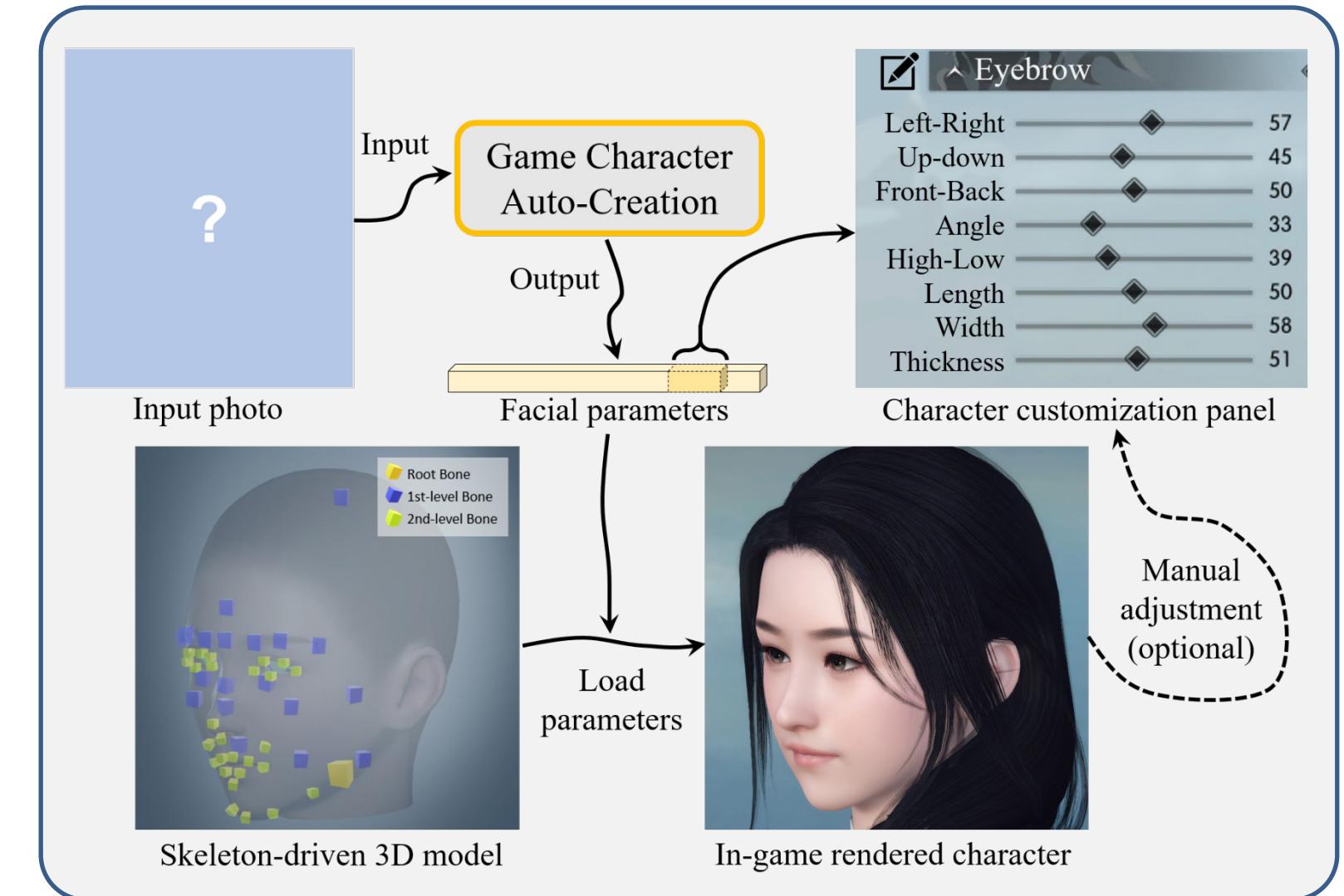
What attributes can be edited?

- Facial shape
- Hairstyles
- More makeups
 - e.g. Eyebrow, Eyeshadow,
 - Pupil color, Blusher, Tattoo,
 - ...



Our goals

- Generate a set of parameters based on the input photo
- These parameters can be recognized by the game
- The created character looks like the input photo
- Players are allowed to further edit the created character



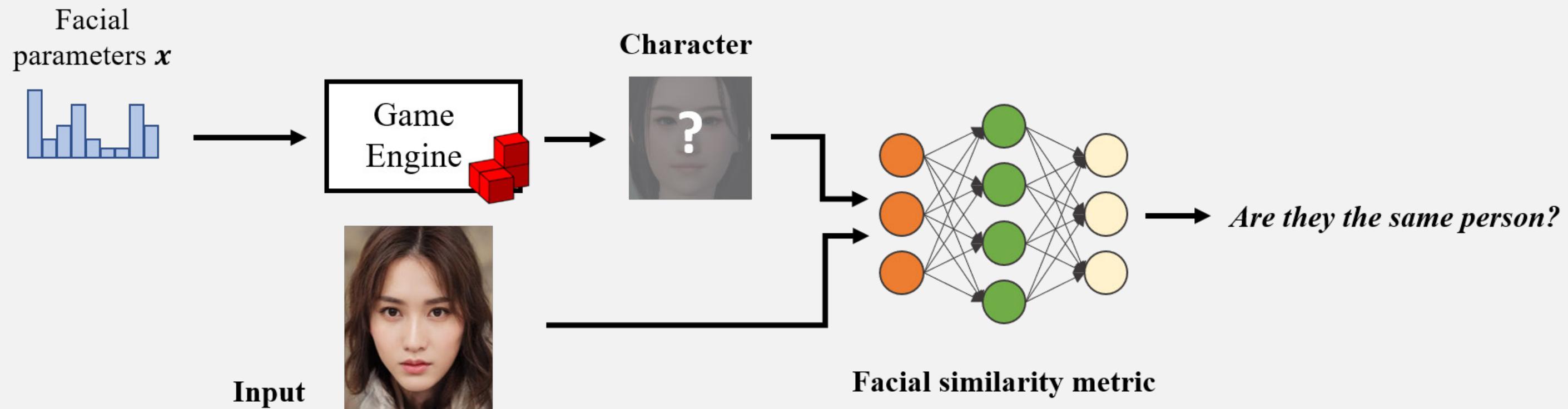


Face-to-Parameter Translation



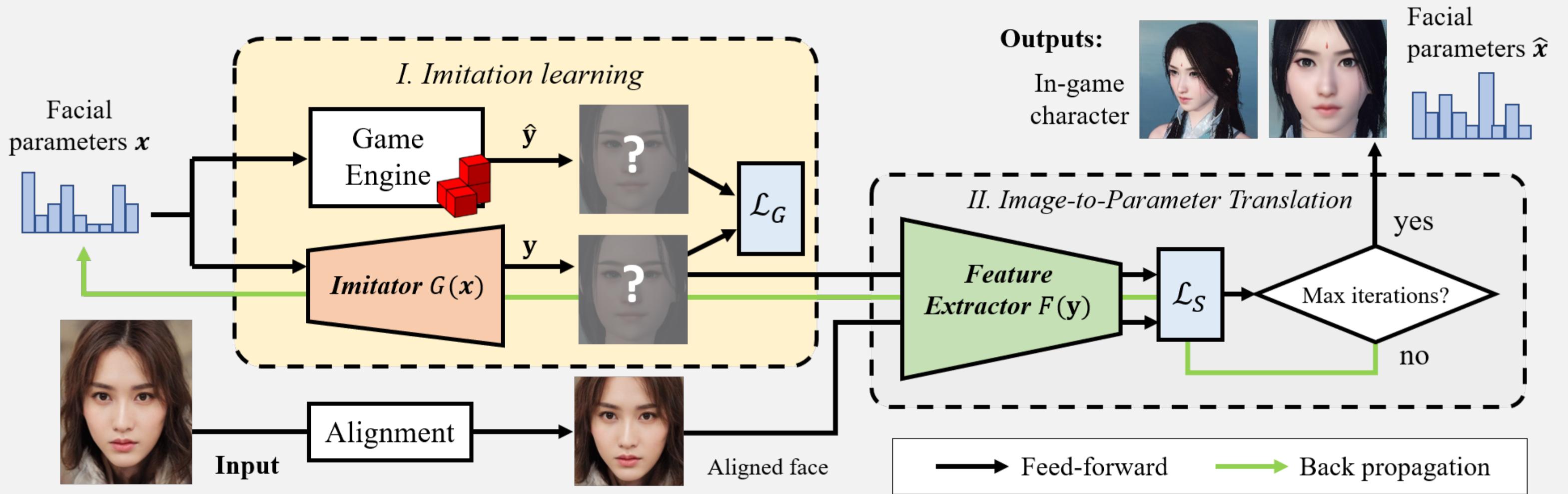
GAME DEVELOPERS CONFERENCE
MARCH 16–20, 2020 | #GDC20

Challenges



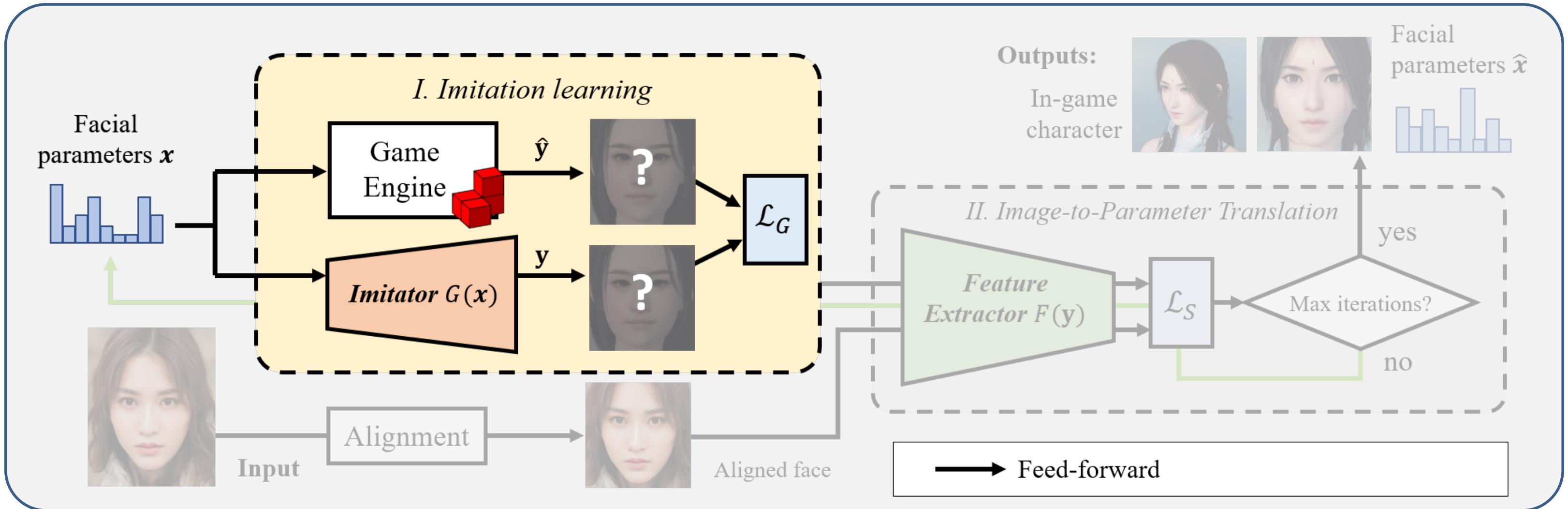
- Game engine is non-differentiable
- Define a metric to measure the similarity between input photo and created character

Methodology



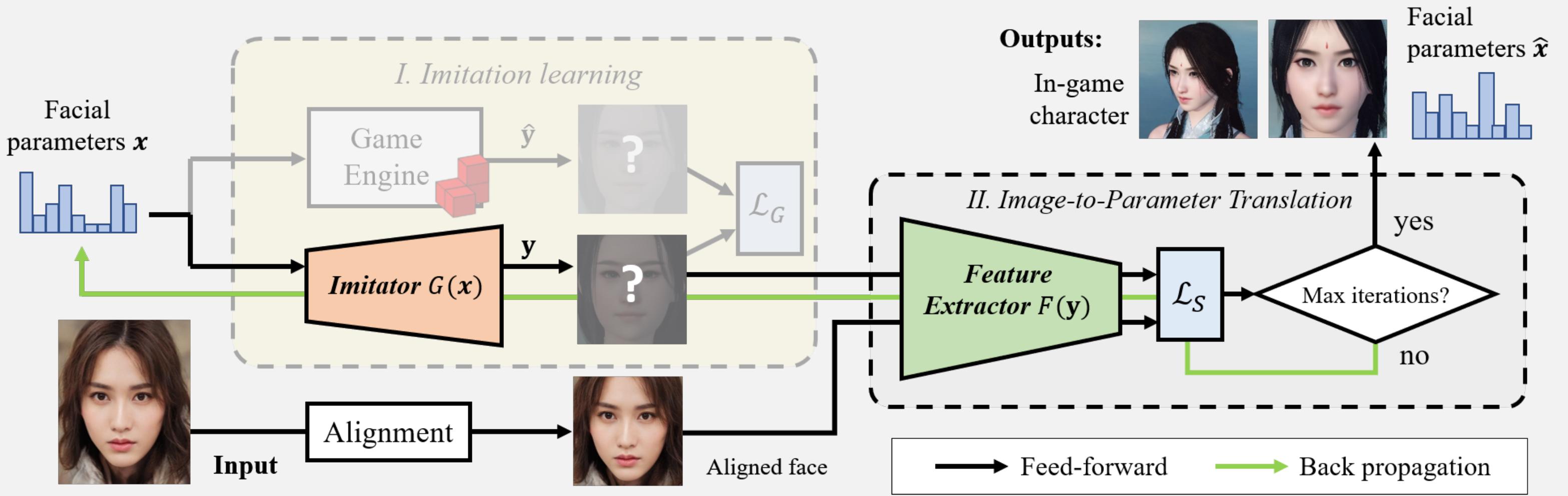
- Phase 1: Train a generative neural network as a renderer
- Phase 2: Define the facial similarity and auto-create characters by the gradient descent method

Methodology



- Phase 1: Train a generative neural network as a renderer
- Phase 2: Define the facial similarity and auto-create characters by the gradient descent method

Methodology



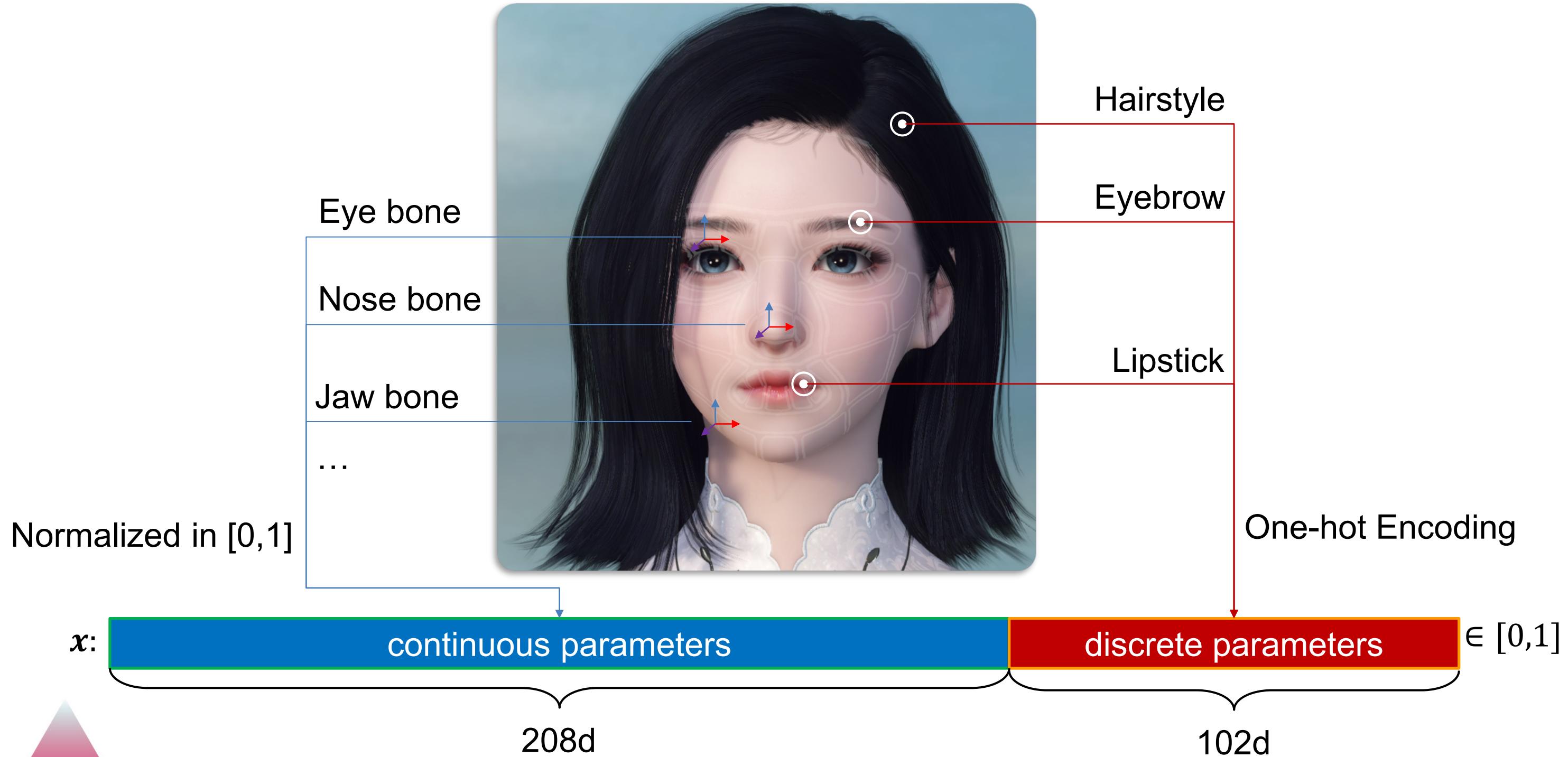
- Phase 1: Train a generative neural network as a renderer
- Phase 2: Define the facial similarity and auto-create characters by the gradient descent method



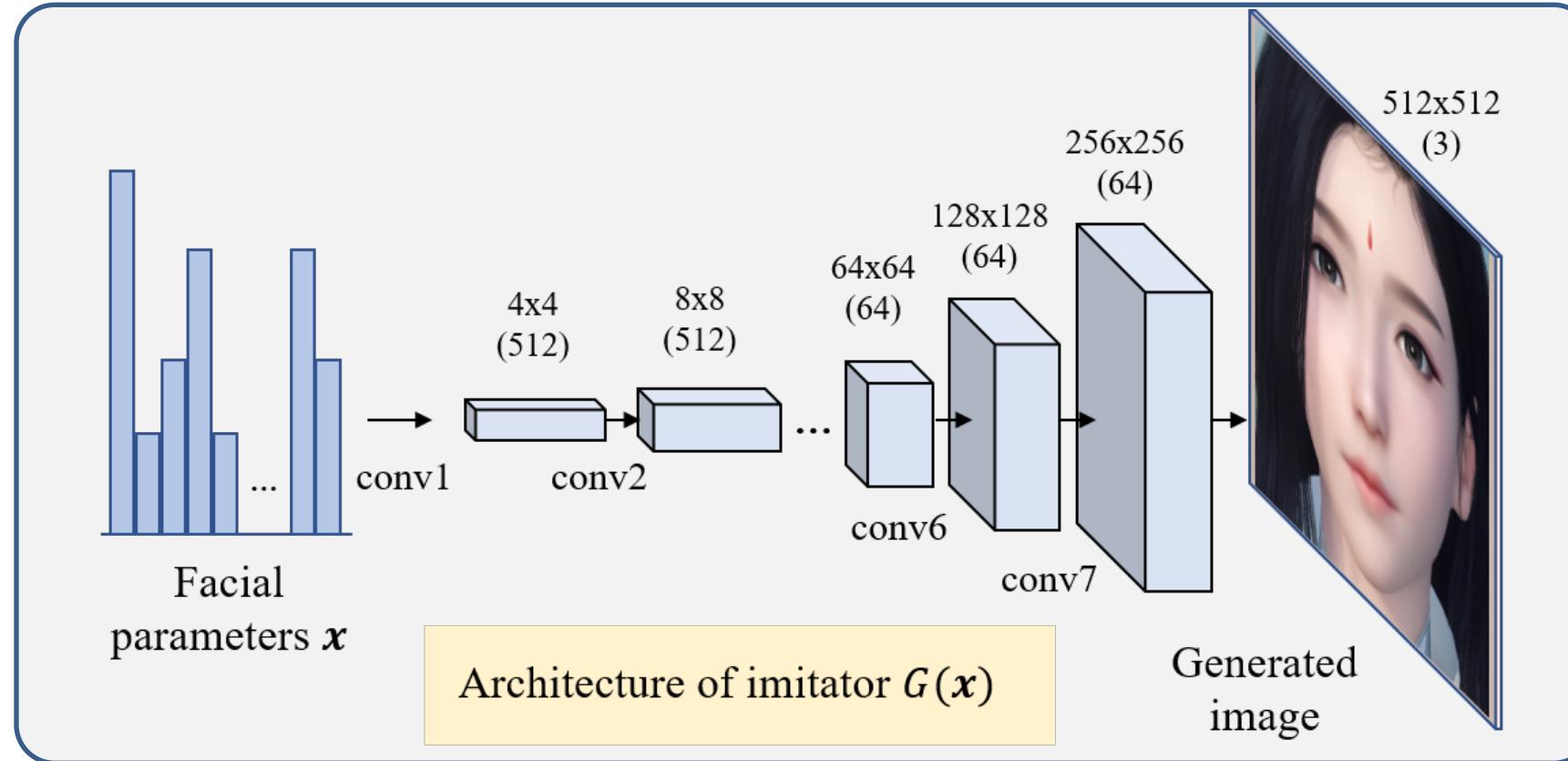
Phase 1: Imitation learning



Facial parameters



Imitation learning

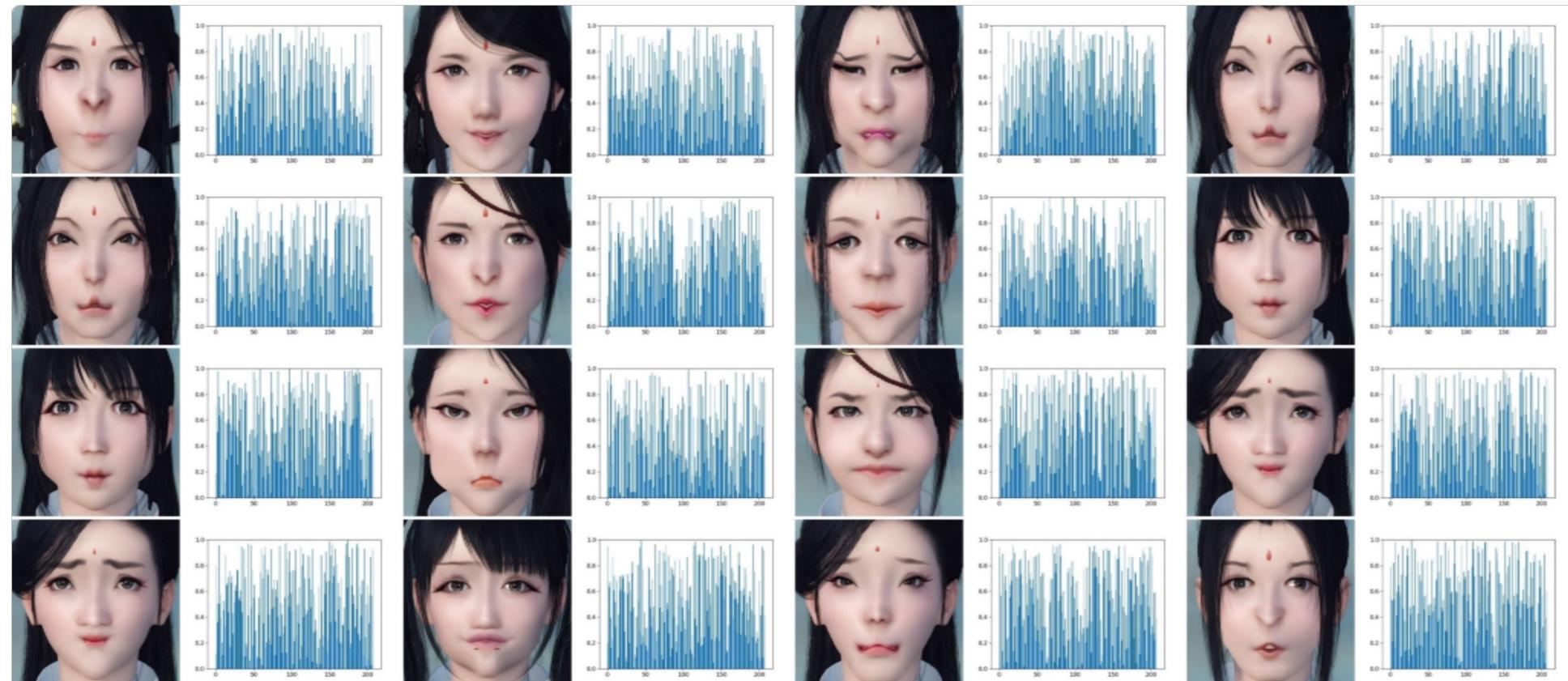


- ## Structure

- Input: facial parameters ($\mathbf{x} \in \mathcal{R}^{208+102}$)
- Output: game character facial image ($G(\mathbf{x}) \in \mathcal{R}^{512 \times 512 \times 3}$)
- 8 transposed convolution layers attached with ReLU and BN layers
- Feature size is doubled after each convolution layer

Imitation learning

- Data
 - Randomly generate facial parameters obeying uniform distribution $x_{continuous} \sim U(0,1)$
 - $x_{discrete} \sim U\{0,1, \dots, N\}$
- Render facial parameters in game (20000 pictures for each kind of character)
- Sample pairs:

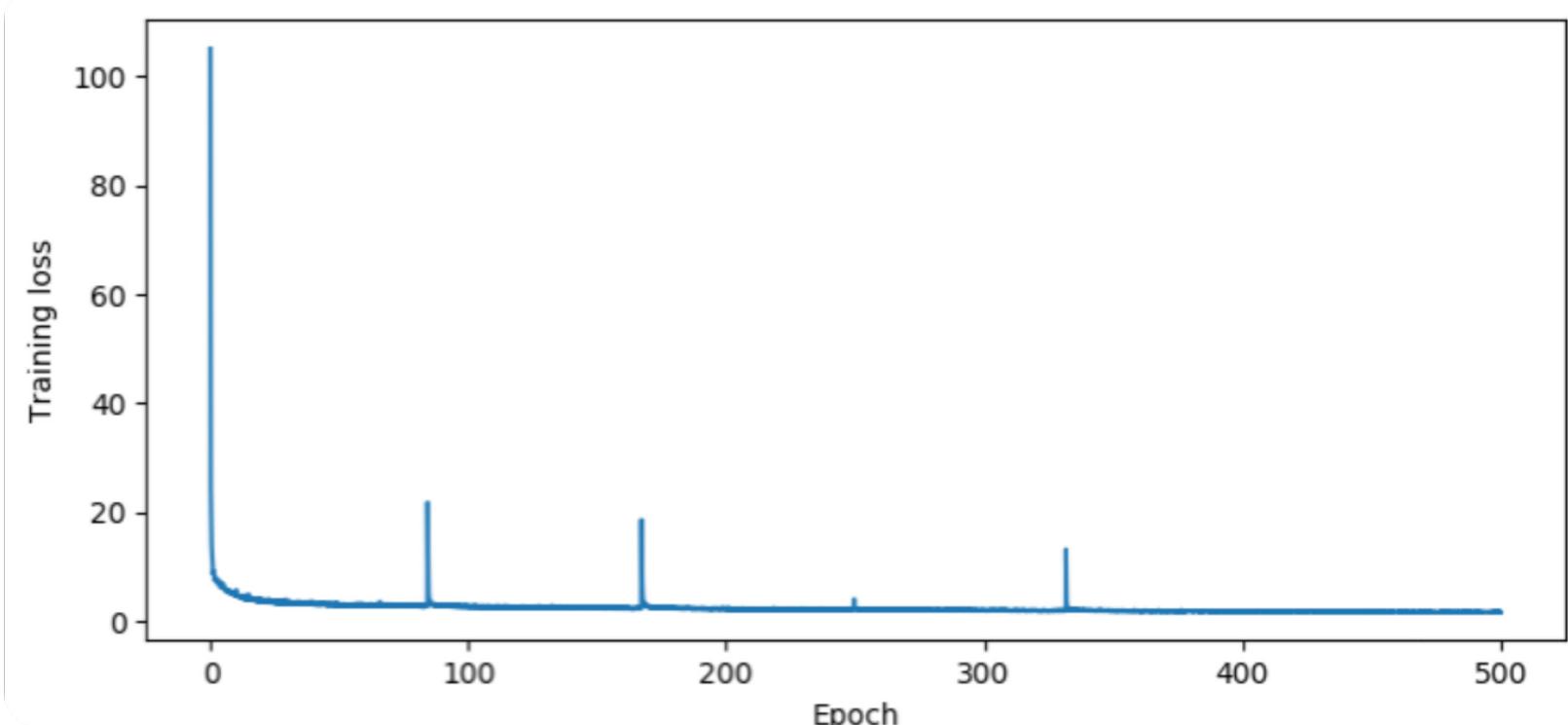


Imitation learning

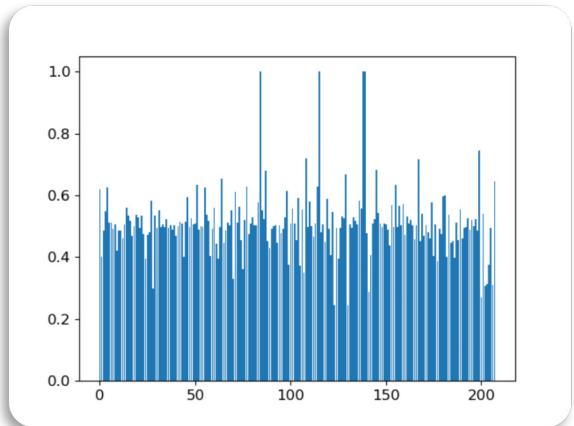
- **Training**
 - Target function:

$$\min_G E_{x \sim u(x)} [\|G(x) - Engine(x)\|_1]$$

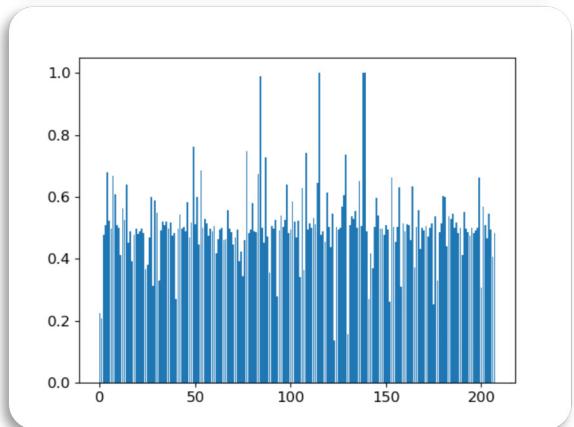
- Optimizer: Stochastic Gradient Descent with momentum
- Configuration: learning rate = 0.01, momentum = 0.9, batch size = 16, maximum epoch = 500, learning rate decay = 10% per 50 epoch



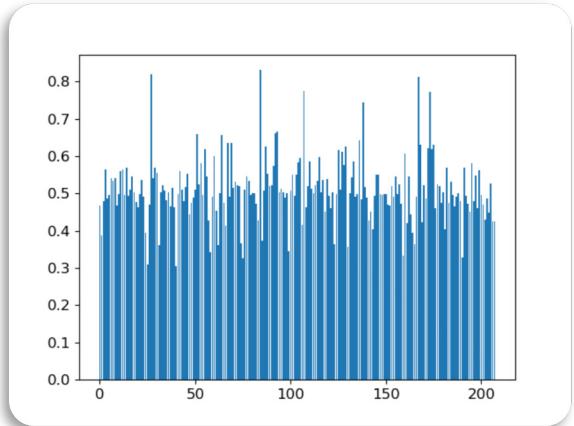
Imitation learning



Imitator
 $G(x)$
→



Imitator
 $G(x)$
→



Imitator
 $G(x)$
→



Input facial parameters x

Generated face

Ground truth

* None of the three pairs are selected from the imitator dataset

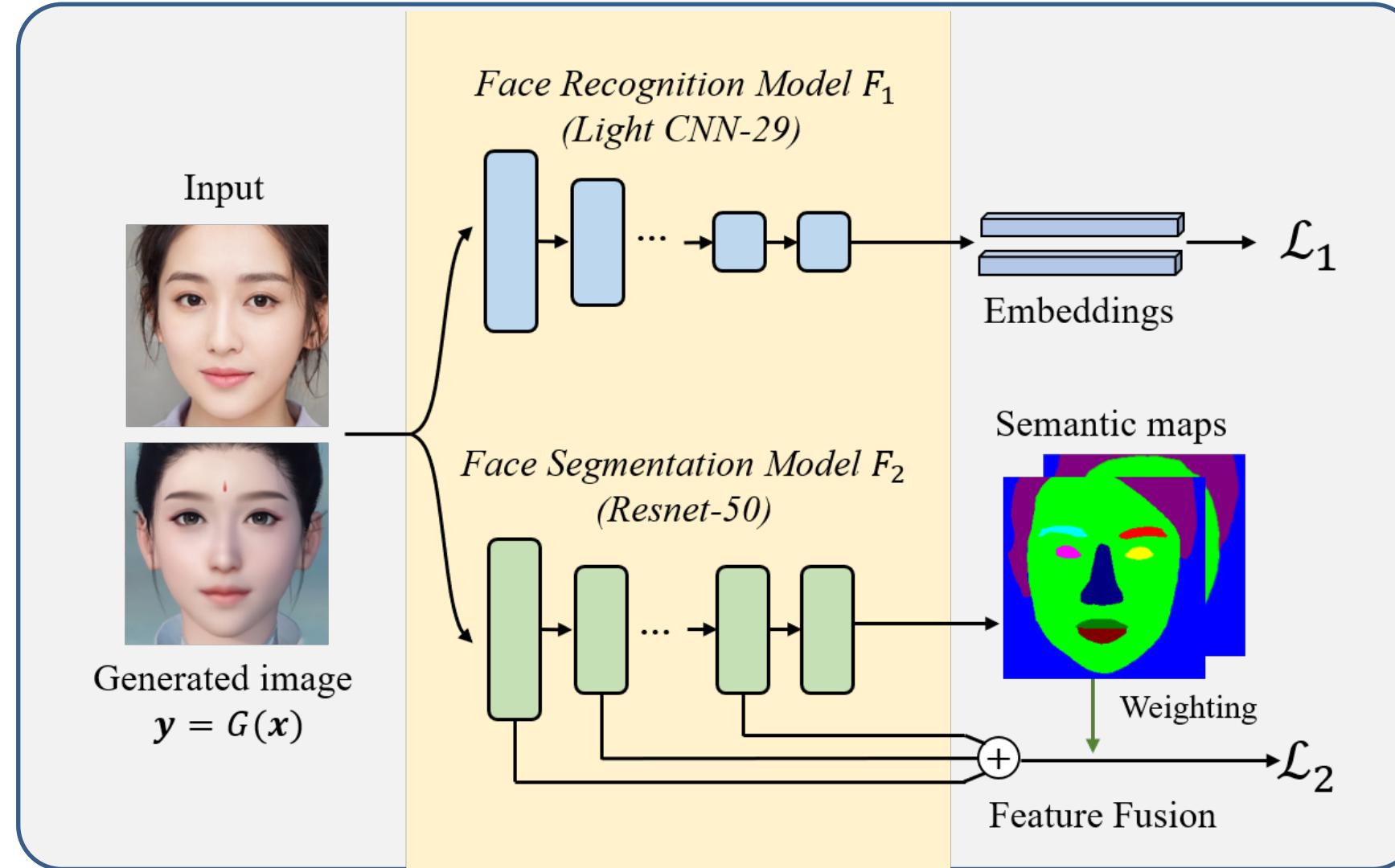


Phase 2: Face-to-Parameter Translation



GAME DEVELOPERS CONFERENCE
MARCH 16–20, 2020 | #GDC20

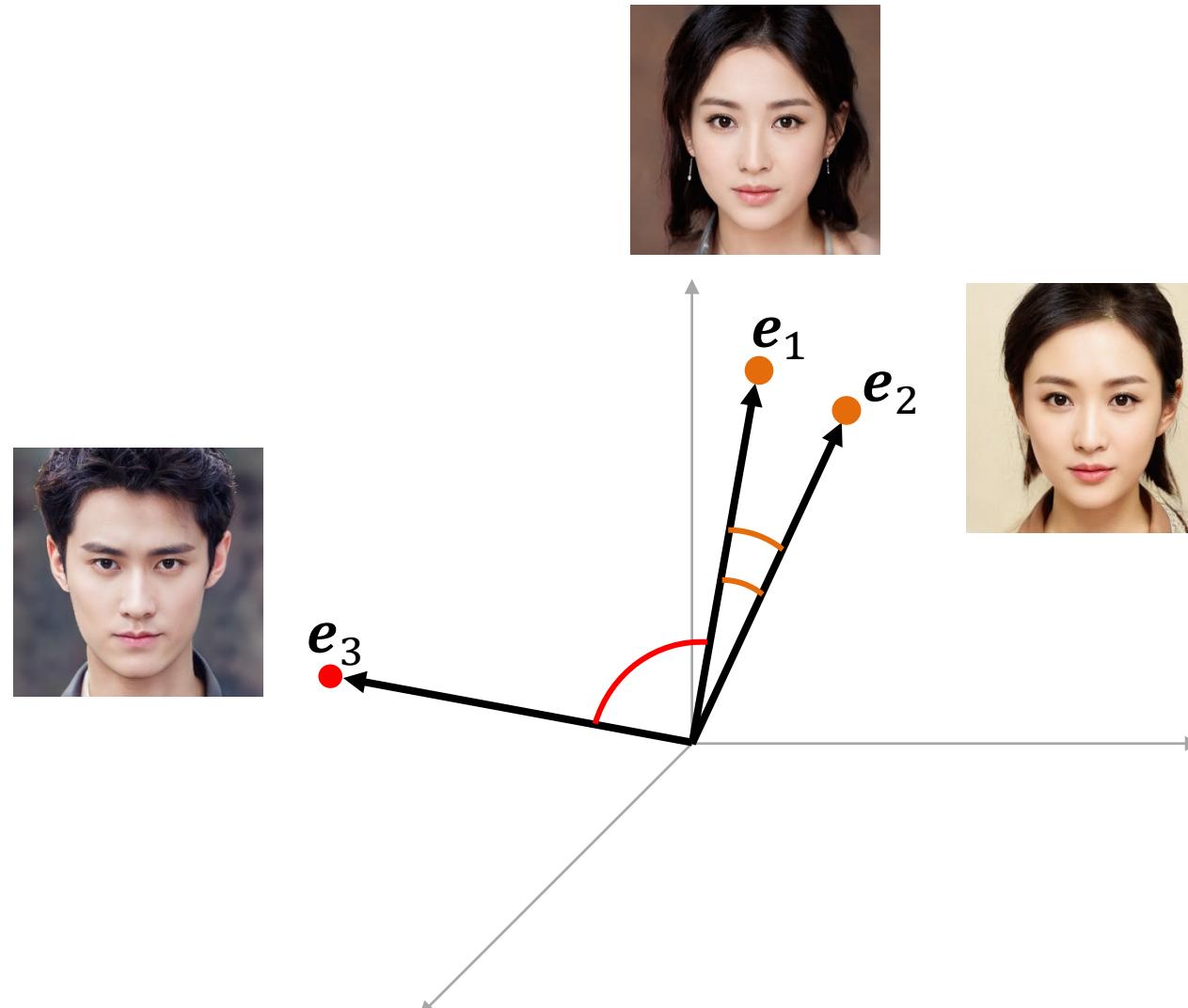
Facial Similarity Measurement



- Use discriminative networks to model human perception
- Input: facial photo (y_r) & generated image ($G(x)$)
- Output: identity loss (\mathcal{L}_1) and content (shape) loss (\mathcal{L}_2)

Face Recognition Model F_1

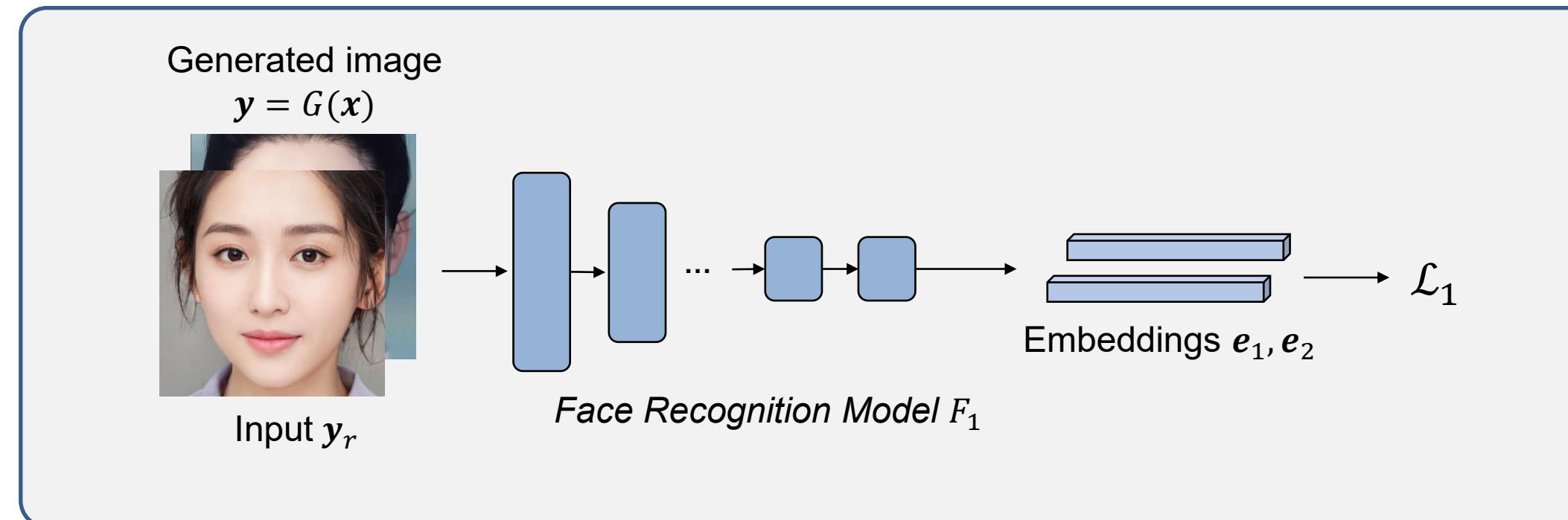
- **Facial embedding**
 - A vector used on identity information representation
 - Photos of the same person share similar embeddings



Face Recognition Model F_1

- **Model**

- Pretrained neural network – Light CNN-29v2*
- Target: Global face information extraction
- Input: 128x128 gray images (input & generated ones)
- Output: 256d facial embeddings
- Fast (~1ms on GTX1080Ti) and accurate (99.43% on LFW**)



* Wu et al., TIFS2018, <https://github.com/AlfredXiangWu/LightCNN>

** Labelled Faces in the Wild (A public face verification dataset)

Face Recognition Model F_1

- Identity loss \mathcal{L}_1

- Definition

$$\mathcal{L}_1(\mathbf{x}) = 1 - \cos \langle F_1(\mathbf{y}_r), F_1(G(\mathbf{x})) \rangle$$

- Cosine distance

$$\cos(\mathbf{e}_1, \mathbf{e}_2) = \frac{\langle \mathbf{e}_1, \mathbf{e}_2 \rangle}{\|\mathbf{e}_1\|_2 \|\mathbf{e}_2\|_2}$$

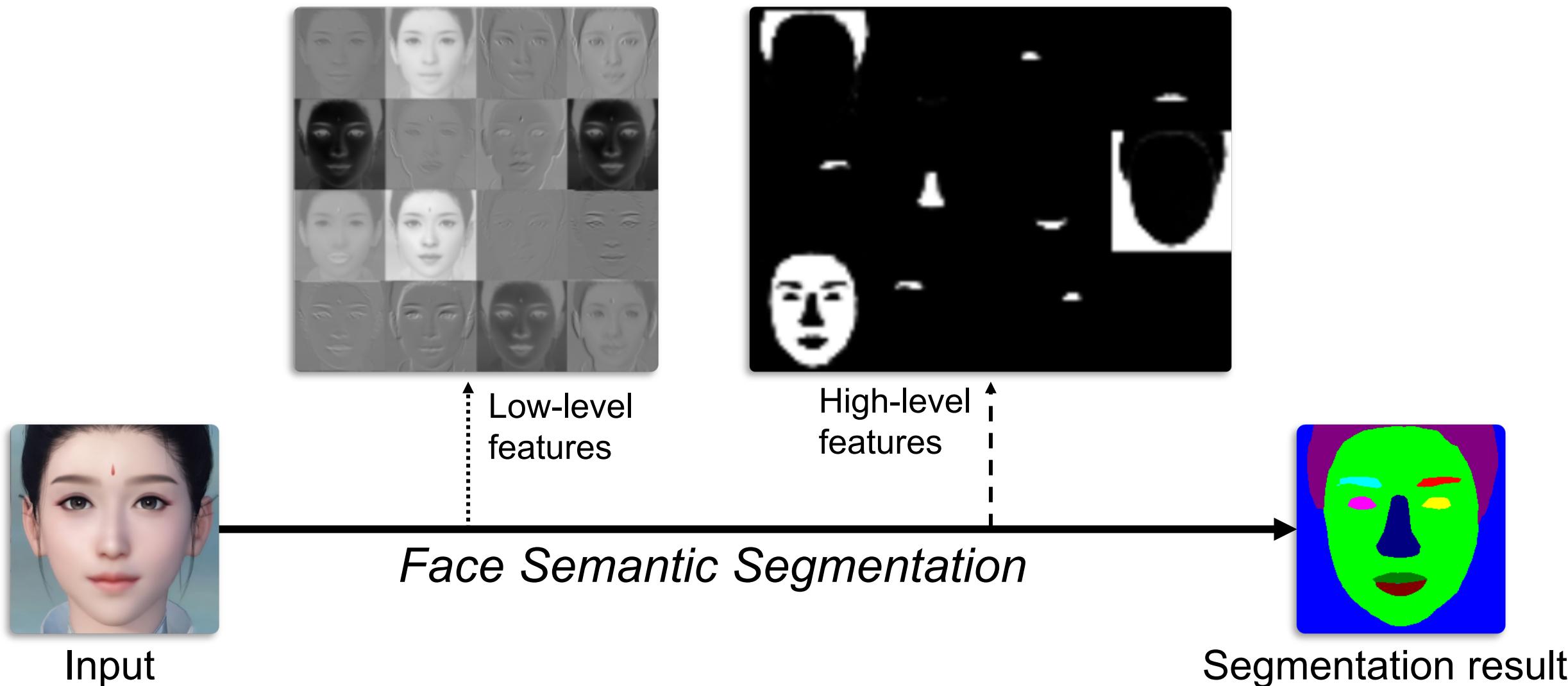
- Optimal object

$$\min_{\mathbf{x}} \mathcal{L}_1(\mathbf{x})$$



Face Semantic Segmentation Model F_2

- **Facial features**
 - A group of 2D features used on facial detail representation



Face Semantic Segmentation Model F_2

- **Model**
 - Target: Local face information extraction
 - Input: 256x256 RGB images
 - Output: facial features and semantic segmentation image
- **Structure**
 - A fully convolution network using Resnet-50 as backbone
 - Replace the last fully connected layer by 1x1 convolution layer
 - Increase output resolution from 1/32 to 1/8

Layer	Component	Configuration	Output Resolution
Segmentation Model	Conv_1	Convolution + BN + ReLU	$64 \times 7 \times 7 / 2$
	MaxPool	MaxPool	$\frac{1}{4} \times \frac{1}{4}$
	Conv_2	3 x Bottleneck	$\frac{1}{8} \times \frac{1}{8}$
	Conv_3	4 x Bottleneck	$\frac{1}{8} \times \frac{1}{8}$
	Conv_4	6 x Bottleneck	$\frac{1}{8} \times \frac{1}{8}$
	Conv_5	3 x Bottleneck	$\frac{1}{8} \times \frac{1}{8}$
	Conv_6	Convolution	$\frac{1}{8} \times \frac{1}{8}$

Resnet-50

Face Semantic Segmentation Model F_2

- **Training Dataset**

- Thousands of human photos with fine label
- 11 classes including eye, nose, mouth, etc.
- Data augmentation by scaling and shifting



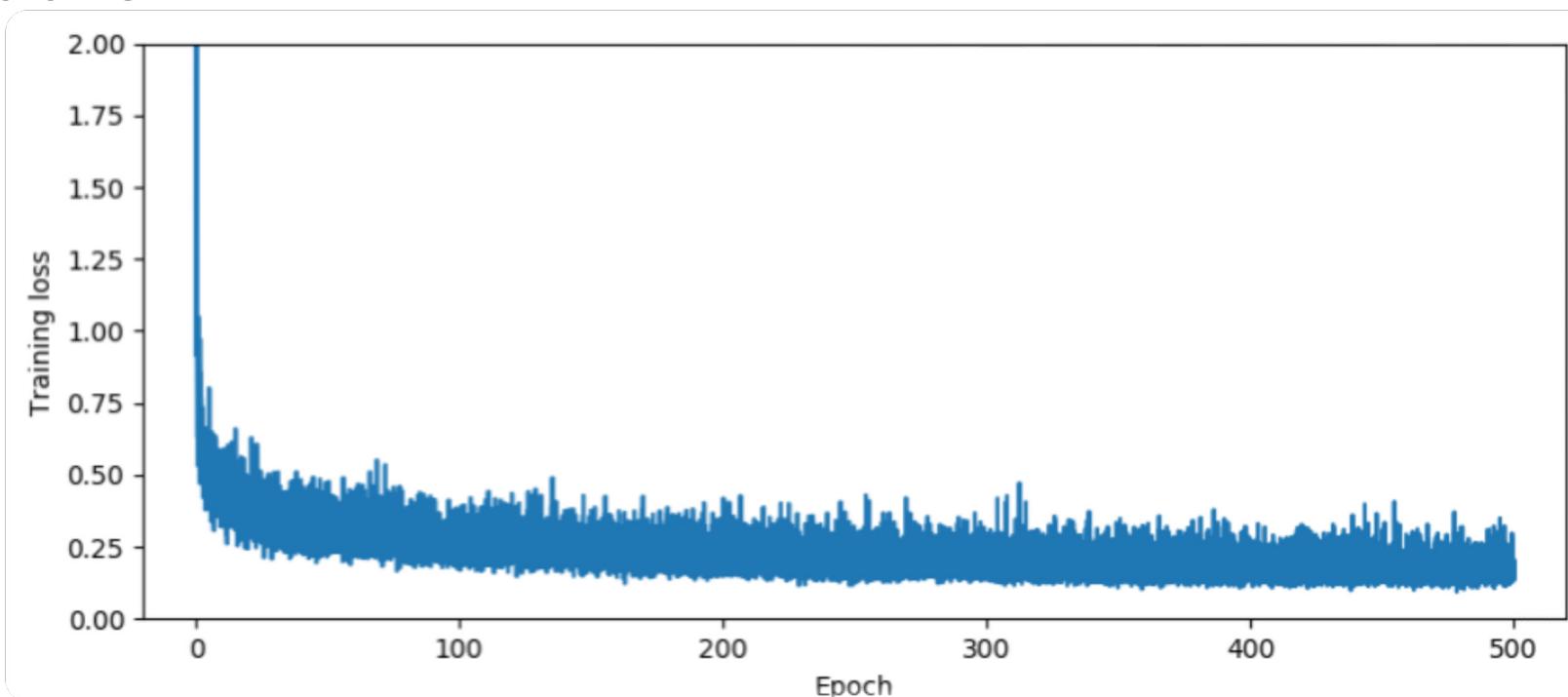
Face Semantic Segmentation Model F_2

- **Training**

- Configuration is similar to the imitator's (except learning rate = 0.001)
- Pixel-wise cross-entropy loss

$$L_{CE} = - \sum_{i=1}^K y_i \log p_i$$

- Training procedure:



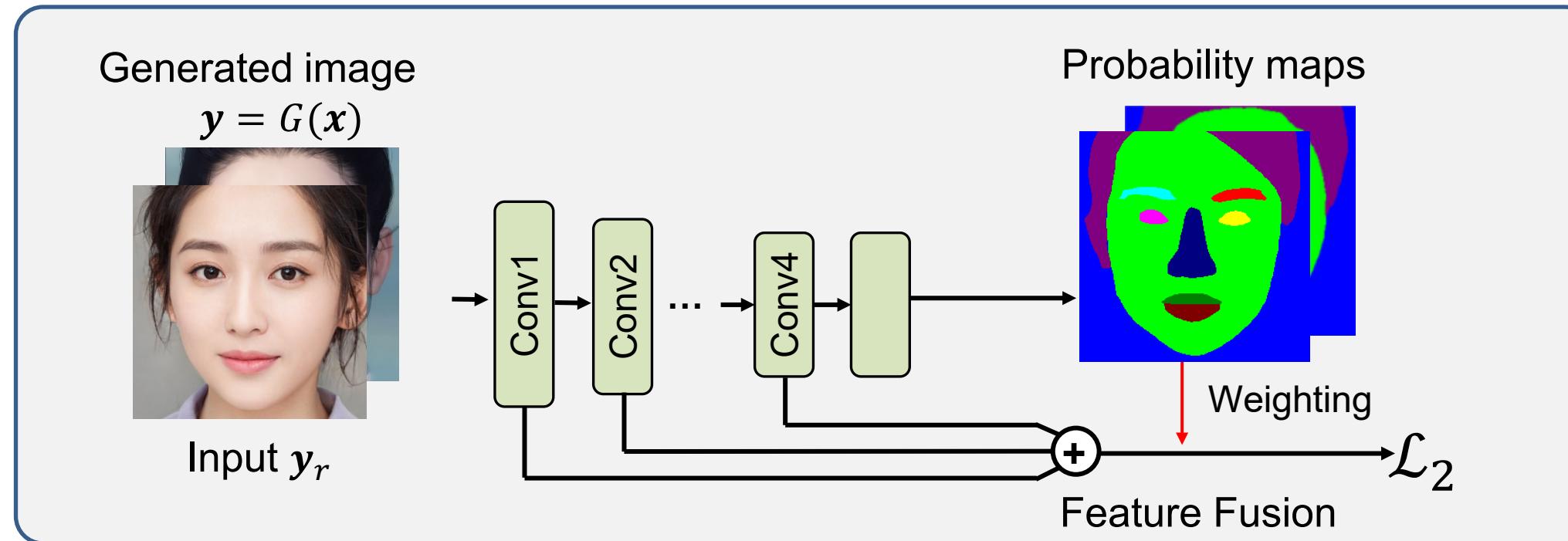
Face Semantic Segmentation Model F_2

- Content loss
 - Definition

$$\mathcal{L}_2(x) = \|F_2(y_r) - F_2(G(x))\|_1$$

- Use semantic probability maps to enhance facial features
- Optimal object

$$\min_x \mathcal{L}_2(x)$$



Face-to-Parameter Translation

- Overall loss function

$$\mathcal{L}_s = \alpha \mathcal{L}_1 + \mathcal{L}_2$$

- Overall optimal problem

$$\min_{\boldsymbol{x}} \mathcal{L}_s(\boldsymbol{x})$$

- Gradient descent method

➤ Parameter update: $\boldsymbol{x}_i = \boldsymbol{x}_{i-1} - \eta \frac{\partial \mathcal{L}_s(\boldsymbol{x}_{i-1})}{\partial \boldsymbol{x}_{i-1}}, \quad i = 1, 2, \dots, i_{max}$

➤ Gradient compute:

$$\frac{\partial \mathcal{L}_s(\boldsymbol{x})}{\partial \boldsymbol{x}} = \alpha \frac{\partial \mathcal{L}_1(\boldsymbol{x})}{\partial \boldsymbol{x}} + \frac{\partial \mathcal{L}_2(\boldsymbol{x})}{\partial \boldsymbol{x}}$$

$$\frac{\partial \mathcal{L}_1(\boldsymbol{x})}{\partial \boldsymbol{x}} = \frac{\partial \mathcal{L}_1(\boldsymbol{x})}{\partial F_1(G(\boldsymbol{x}))} \frac{\partial F_1(G(\boldsymbol{x}))}{\partial G(\boldsymbol{x})} \frac{\partial G(\boldsymbol{x})}{\partial \boldsymbol{x}}$$

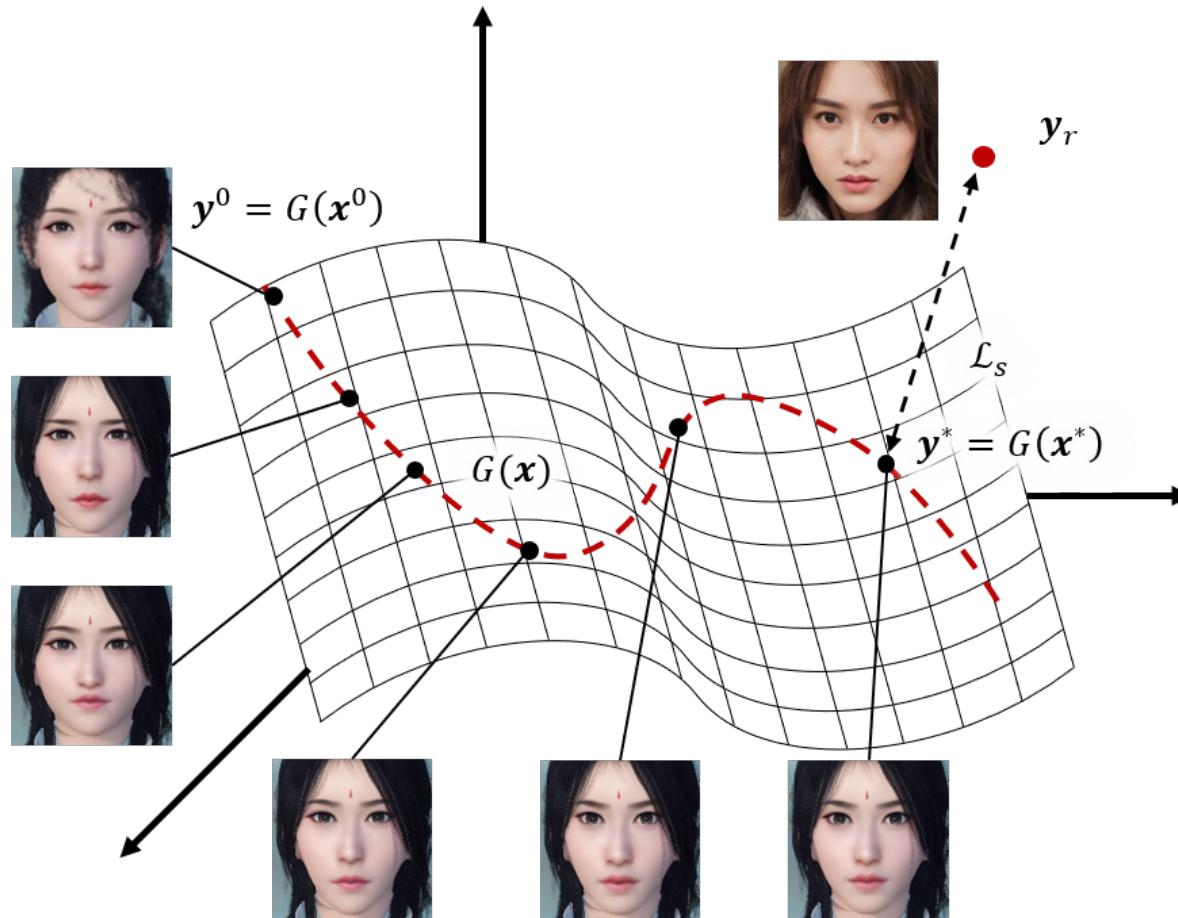
$$\frac{\partial \mathcal{L}_2(\boldsymbol{x})}{\partial \boldsymbol{x}} = \frac{\partial \mathcal{L}_2(\boldsymbol{x})}{\partial F_2(G(\boldsymbol{x}))} \frac{\partial F_2(G(\boldsymbol{x}))}{\partial G(\boldsymbol{x})} \frac{\partial G(\boldsymbol{x})}{\partial \boldsymbol{x}}$$



Manifold searching

- Geometrically, the gradient descent procedure can be viewed as a searching on the game face manifold:

$$\Rightarrow \min_{\mathbf{x}} \alpha(1 - \cos \langle F_1(\mathbf{y}_r), F_1(\mathbf{y}) \rangle) + \|F_2(\mathbf{y}_r) - F_2(\mathbf{y})\|_1 \\ s.t. \quad \mathbf{y} = G(\mathbf{x})$$



Character Auto-Creation

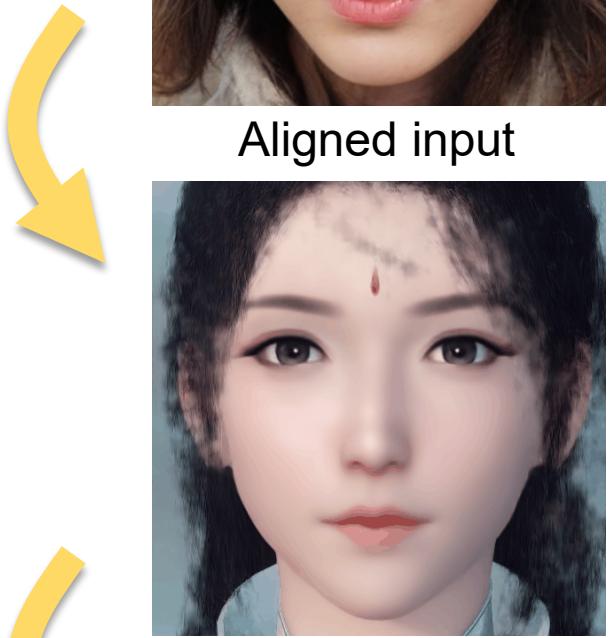
- Step 1: Align the user-uploaded photo based on the default game face
- Step 2: Solve the following optimization problem

$$\min_x \mathcal{L}_s(x)$$

- Step 2.1: Freeze all networks ($G F_1 F_2$), initialize $x = x_0$, set learning rate $\eta = 10$ and max iterations $i_{max} = 10$
- Step 2.2: Use gradient descent method to update x
- Step 2.3: When reaching the max iterations, output optimized facial parameters x^*
- Step 3: Write x^* into the game and obtain the character



Aligned input



Optimizing



In-game character

Player editable

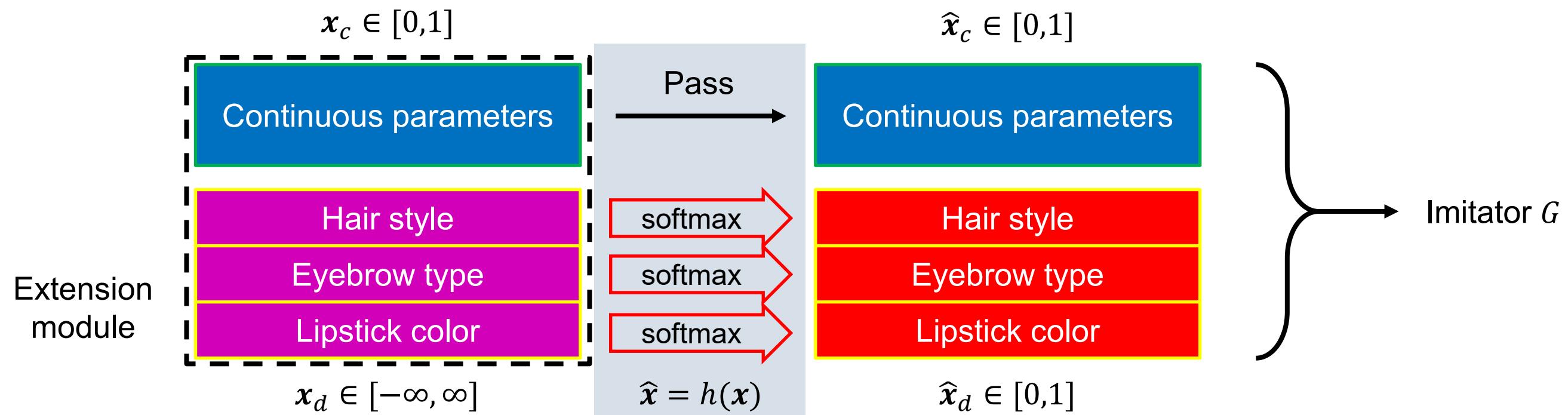


Changes:
Eye shadow
Lipstick color



The solution of discrete parameter optimization

- Discrete parameters are hard to optimize directly
- SoftMax function ($\text{softmax}(\mathbf{x}) = \frac{e^{\mathbf{x}}}{\sum e^{\mathbf{x}}}$) is differentiable and can approximate one-hot vector by generating a probability vector

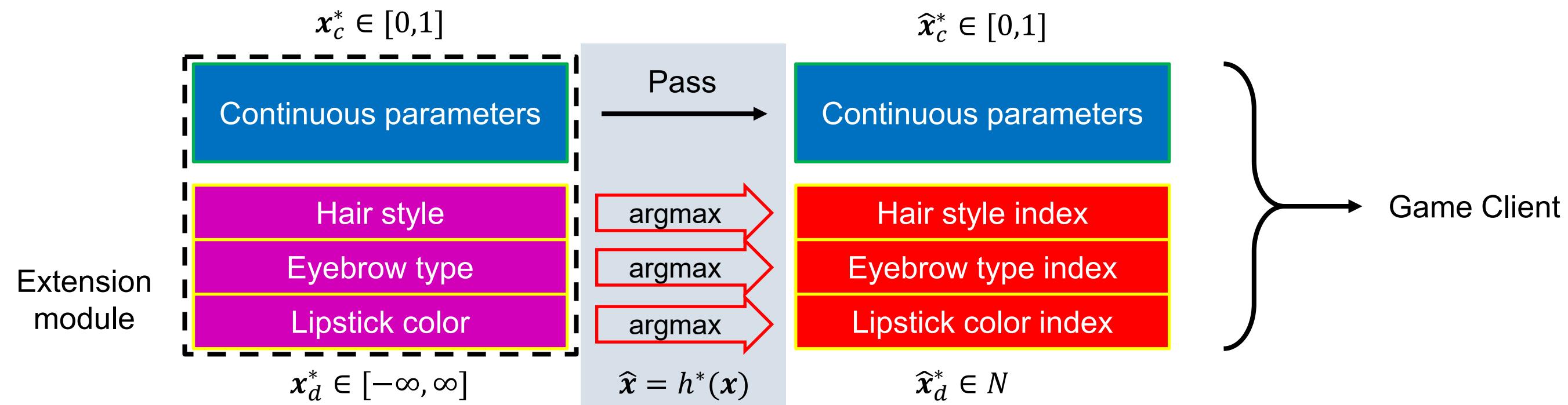


- New optimal problem

$$\min_{\mathbf{x}} \mathcal{L}_s(h(\mathbf{x}))$$

The solution of discrete parameter optimization

- When obtaining the optimized facial parameters \mathbf{x}^* , softmax function will be replaced by argmax function to get the index of makeups



Note: If $y = \text{argmax}(\mathbf{x})$, then y is the index of the maximum element in the vector \mathbf{x}

Environment and Performance

Action	Software	Hardware	Speed
Imitator training dataset	Windows 10, Game client	I7, RAM 16GB, GTX1060	20,000 pictures / 10 hours
Imitator training		E5, RAM 128GB, GTX1080Ti x4	1 day
Segmentation network training	Linux, PyTorch	E5, RAM 32GB, GTX1080Ti x1	1 day
Character Auto-Creation		E5, RAM 16GB, GTX1080Ti x1	1 fps

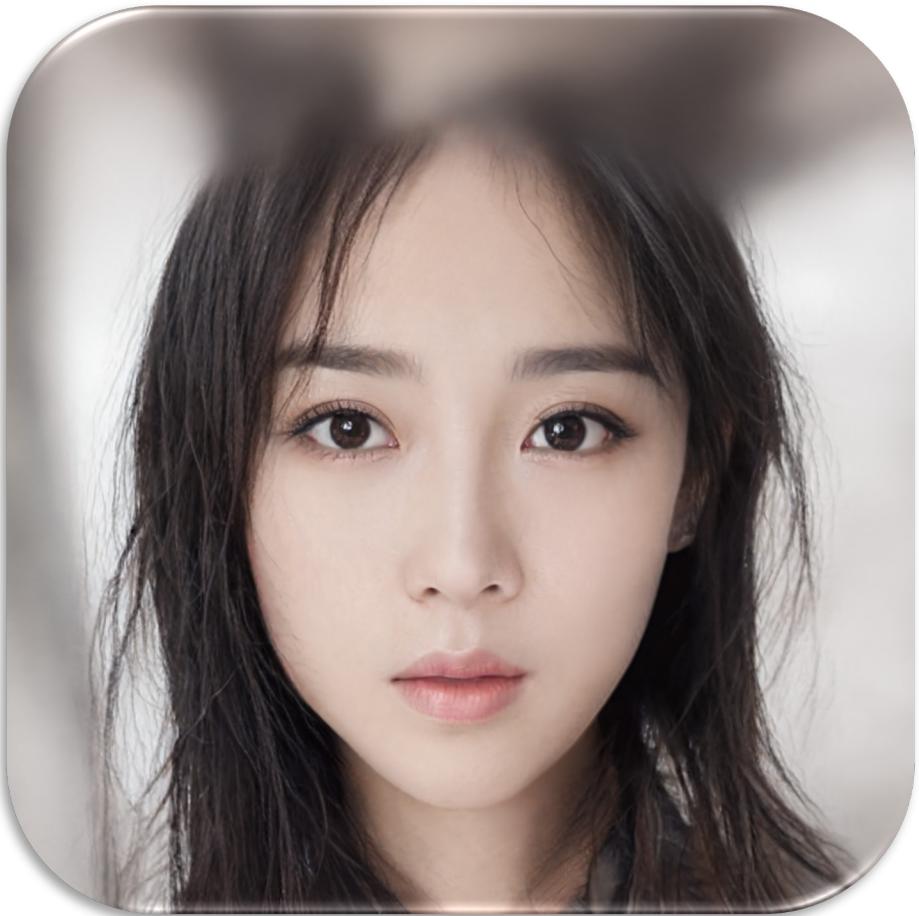


Auto-creation results



GAME DEVELOPERS CONFERENCE
MARCH 16–20, 2020 | #GDC20

“Justice”



GAME DEVELOPERS CONFERENCE
MARCH 16–20, 2020 | #GDC20



“Justice”



GAME DEVELOPERS CONFERENCE
MARCH 16–20, 2020 | #GDC20



“Justice”



GAME DEVELOPERS CONFERENCE
MARCH 16–20, 2020 | #GDC20



“Revelation Mobile”



GAME DEVELOPERS CONFERENCE
MARCH 16–20, 2020 | #GDC20



“Revelation Mobile”



GDC

GAME DEVELOPERS CONFERENCE
MARCH 16–20, 2020 | #GDC20





Discussion



GAME DEVELOPERS CONFERENCE
MARCH 16–20, 2020 | #GDC20

Advantages

- No labeled character data is required
- The generative network can easily fit the behavior of game engine regardless of the renderer type and 3D model structure
- Deep learning based features can represent human perception on face similarity evaluation

Limitations

- Sensitive to the pose, obstacle, etc.
- Discrete parameters are not well solved
- Human-like 3D model is required

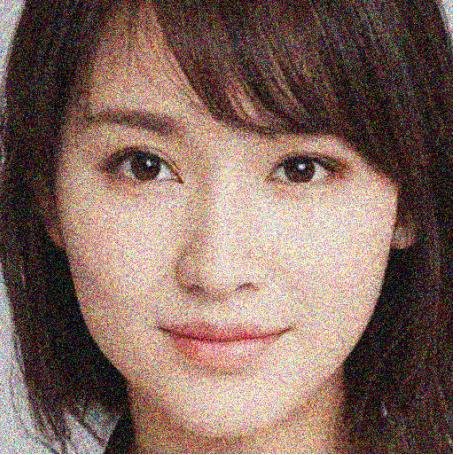


Robustness experiment

Original



Noise



Illumination



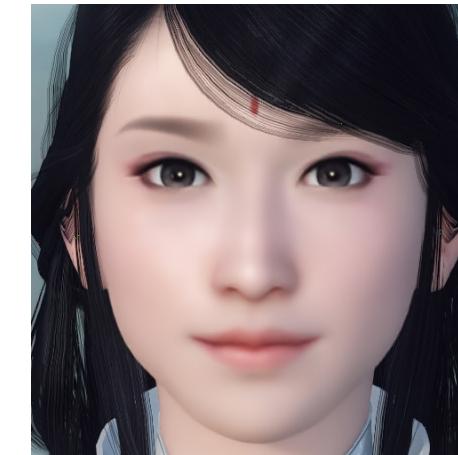
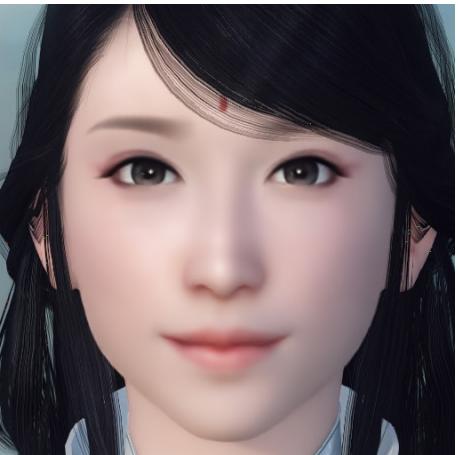
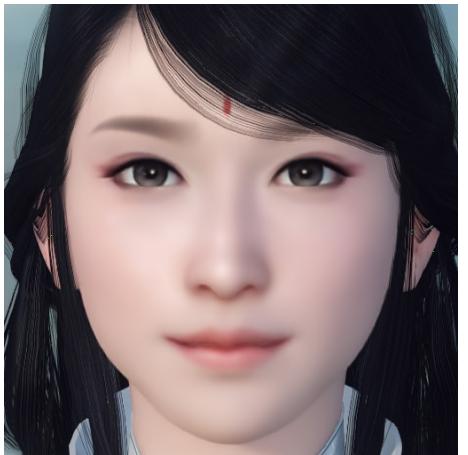
Color



Blur



Pose



(Failure case)



Thanks!



GAME DEVELOPERS CONFERENCE
MARCH 16–20, 2020 | #GDC20