SIGGRAPH 2022
VANCOUVER+ 8-11 AUG

NANYANG TECHNOLOGICAL UNIVERSITY SINGAPORE

S-LAB FOR ADVANCED INTELLIGENCE

商汤 sensetime

上海人工智能实验室 Shanghai Artificial Intelligence Laboratory

THE PREMIER CONFERENCE & EXHIBITION ON COMPUTER GRAPHICS & INTERACTIVE TECHNIQUES

Fangzhou Hong[1*]   Mingyuan Zhang[1*]   Liang Pan[1]   Zhongang Cai[1,2,3]   Lei Yang[2]   Ziwei Liu[1+]

[1]S–Lab Nanyang Technological University   [2]SenseTime Research   [3]Shanghai AI Laboratory
*Equal Contribution   +Corresponding Author

# AvatarCLIP
# ZERO-SHOT TEXT-DRIVEN GENERATION AND ANIMATION OF 3D AVATARS

DALL·E [1]

DALL·E 2 [2]

Imagen [3]

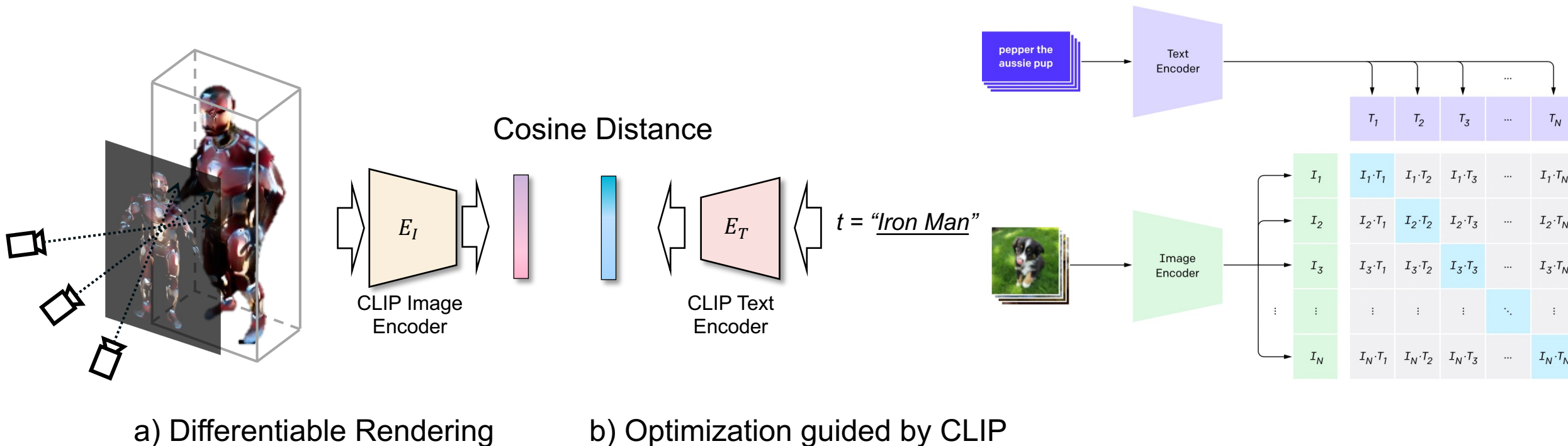[1] https://openai.com/blog/dall-e/      [2] https://openai.com/dall-e-2/
[3] https://imagen.research.google

# TEXT-DRIVEN 3D GENERATION

## CLIP + DIFFERENTIABLE RENDERING

Cosine Distance

$E_I$

CLIP Image
Encoder

$E_T$

CLIP Text
Encoder

$t$ = "*Iron Man*"

a) Differentiable Rendering

b) Optimization guided by CLIP

pepper the
aussie pup

Text
Encoder

Image
Encoder

# TEXT-DRIVEN 3D GENERATION

## CLIP + DIFFERENTIABLE RENDERING

Dream Field [1]



Text2Mesh [2]

[1] https://ajayj.com/dreamfields        [2] https://threedle.github.io/text2mesh/

# WHAT ABOUT TEXT-DRIVEN AVATAR GENERATION => NOW WE HAVE AVATARCLIP

*I want to generate a tall and fat Iron Man that is running.*

*I would like to generate a skinny ninja that is raising arms.*

*I want to generate a tall and skinny female soldier that is arguing.*

*I want to generate an overweight sumo wrestler that is sitting.*

# AVATARCLIP: HOW IT WORKS
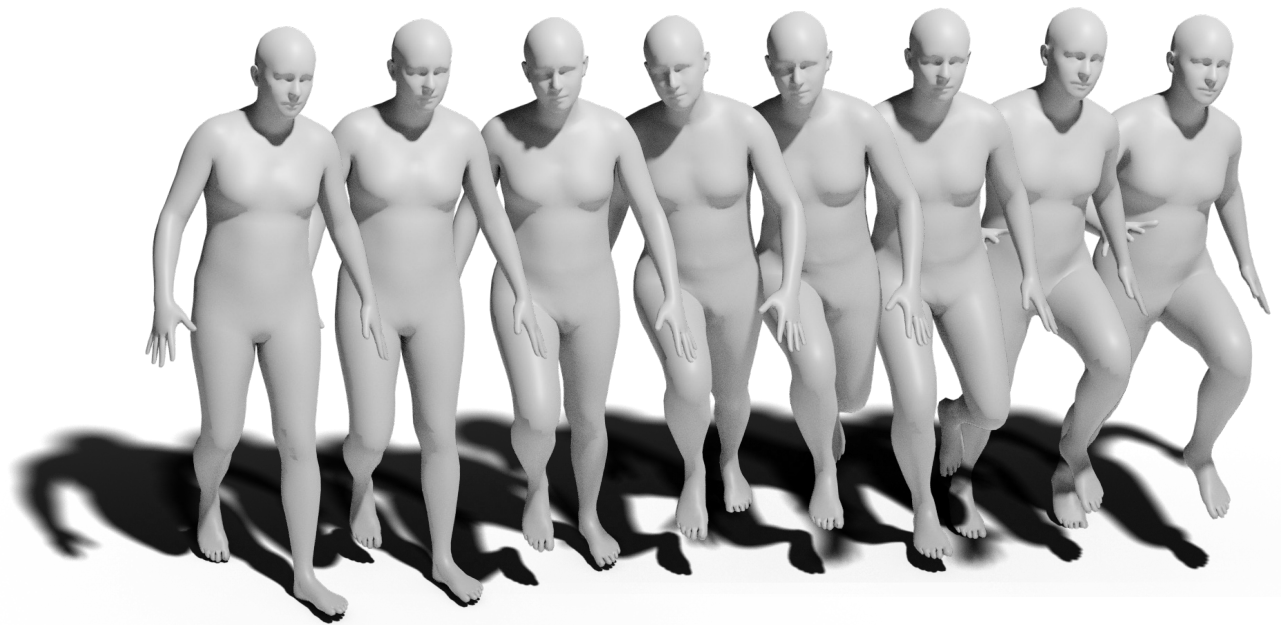
## A) STATIC AVATAR GENERATION

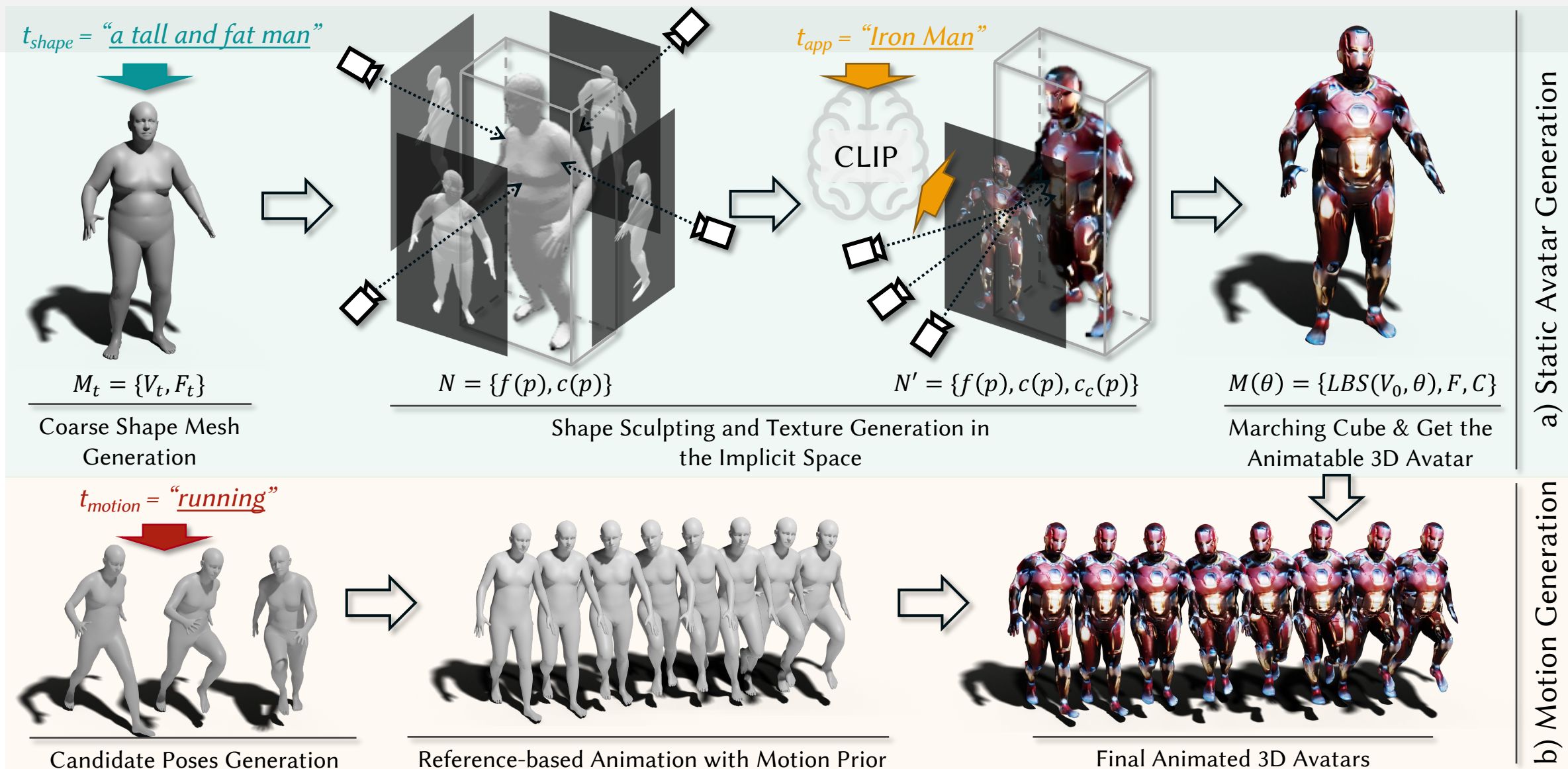Shape Description: *"a tall and fat man"*

Appearance Description: *"Iron Man"*

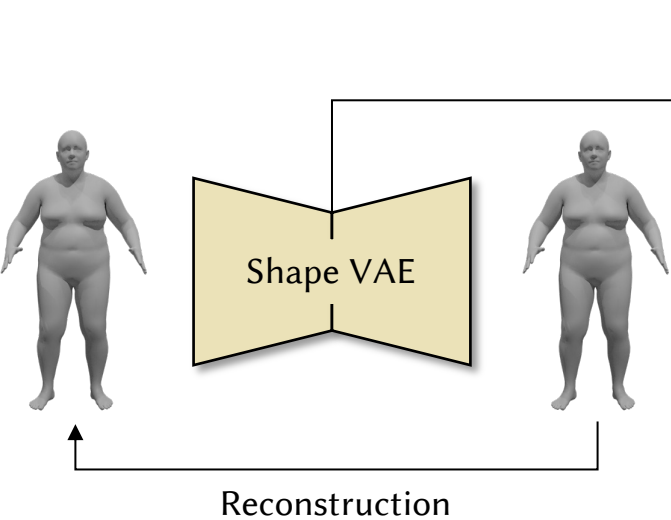## B) MOTION GENERATION

Motion Description: *"running"*
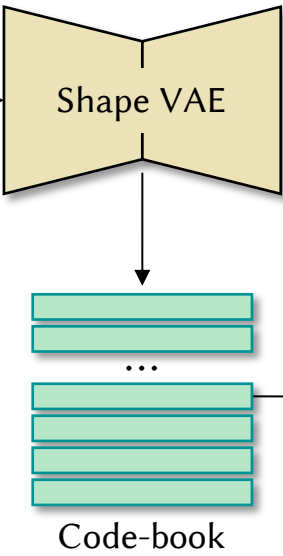
# AVATARCLIP: DETAILED PIPELINE

$t_{shape}$ = "a tall and fat man"

$t_{app}$ = "Iron Man"

CLIP

$t_{motion}$ = "running"

$M_t = \{V_t, F_t\}$

$N = \{f(p), c(p)\}$

$N' = \{f(p), c(p), c_c(p)\}$

$M(\theta) = \{LBS(V_0, \theta), F, C\}$

Coarse Shape Mesh Generation

Shape Sculpting and Texture Generation in the Implicit Space

Marching Cube & Get the Animatable 3D Avatar

a) Static Avatar Generation

Candidate Poses Generation

Reference-based Animation with Motion Prior

Final Animated 3D Avatars

b) Motion Generation

1) a tall man  2) a very skinny man  3) an overweight man

(i)  (ii)  Ours  (i)  (ii)  Ours  (i)  (ii)  Ours

a) Comparison with Baseline Methods

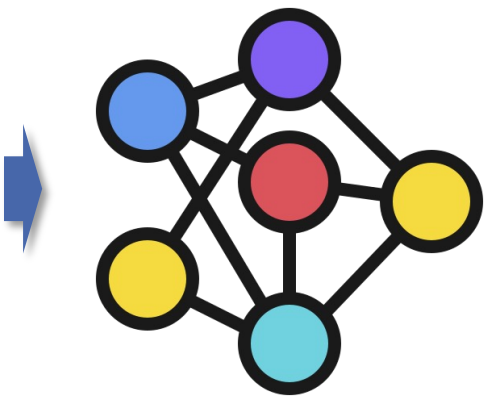a fat man  a short man  a strong man  a short woman  a tall and skinny woman  an overweight woman  an overweight and tall woman

b) More Results

(i) Direct optimization on SMPL parameter beta

(ii) Direct optimization on shape VAE latent code

Mesh

NeuS

$$N = \{f(p), c(p)\}$$

$$C(o, v) = \int_0^\infty w(t) c(p(t)) \, dt$$

$$I_t^{i'}$$

Implicit Function

a) Rendering the Implicit 3D Avatar $N' = \{f(p), c(p), c_c(p)\}$

b) Optimization
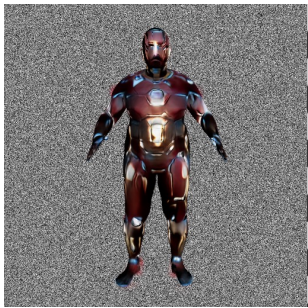
Examples of Intermediate Results
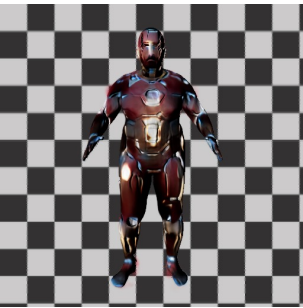
## A) RANDOM BACKGROUND AUGMENTATION



1) Black    2) White

3) Gaussian Noise    4) Chess Board

## B) SEMANTIC-AWARE PROMPT AUGMENTATION
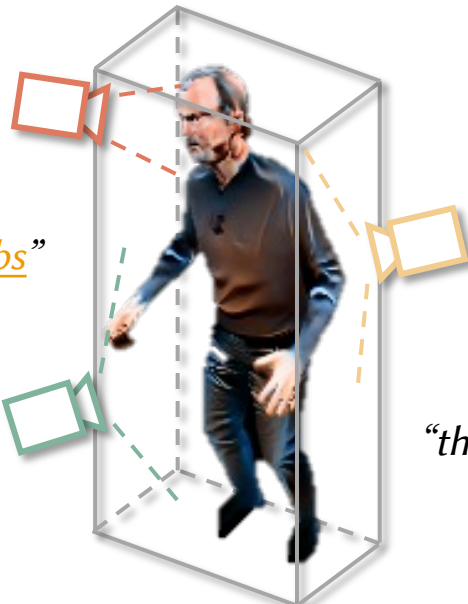


*"the face of <u>Steve Jobs</u>"*

*"the back of <u>Steve Jobs</u>"*

*"<u>Steve Jobs</u>"*

Implicit 3D Avatar $N' = \{f(p), c(p), c_c(p)\}$

Donald Trump

Gentleman

Queen Elizabeth

Ablation 1

Baseline

Ablation 2

+ Background
Augmentation

Ablation 3

+ Texture-less
Renderings

Ablation 4

+ Shading on
Textured Renderings

Full Model

+ Semantic-Aware
Prompt Augmentation

# AVATAR GENERATION RESULTS

Elvis Presley

Freddie Mercury

Drake

Ellen DeGeneres

Karl Lagerfeld

Simon Cowell

1. Superman
2. the face of Bill Gates

1. Iron Man
2. the face of Steve Jobs

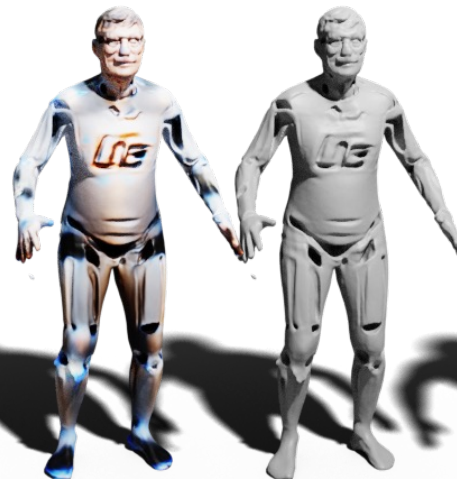Steve Jobs in White Shirt

Man in Jeans

Man in White Shirt

Alien Bill Gates

Bill Gates Wearing Batman Suit
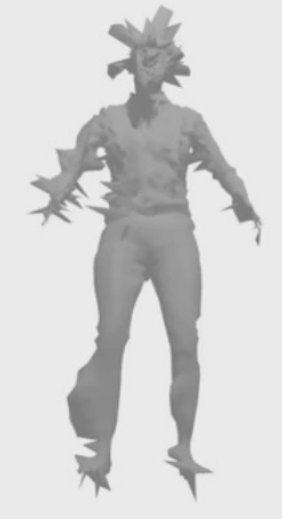
Robot Bill Gates

Zombie Steve Jobs

Zombie Iron Man

# COMPARISON WITH BASELINE METHODS OF AVATAR GENERATION
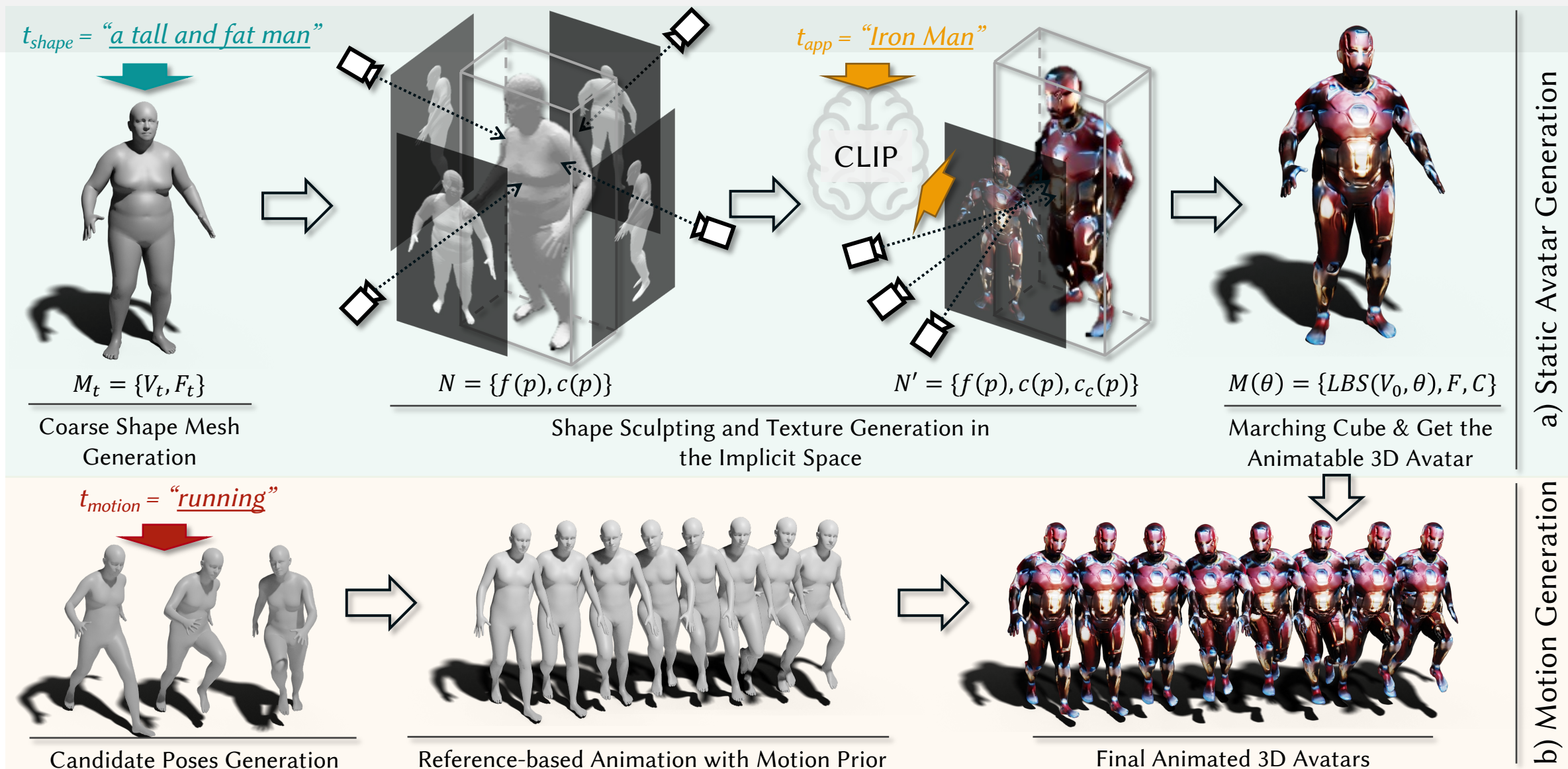
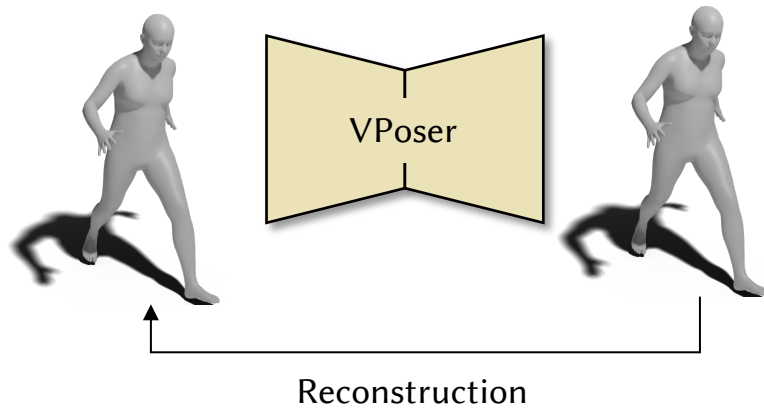Abraham Lincoln

Elvis Presley

Ours | Dream Field (Adapted) | Text2Mesh

$t_{shape}$ = "a tall and fat man"
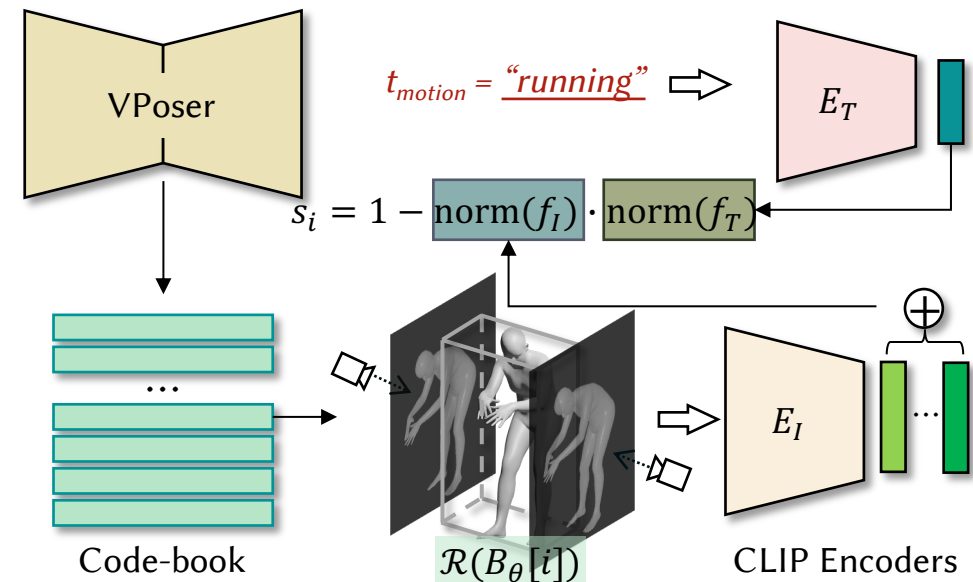
$t_{app}$ = "Iron Man"

CLIP

$M_t = \{V_t, F_t\}$

$N = \{f(p), c(p)\}$

$N' = \{f(p), c(p), c_c(p)\}$

$M(\theta) = \{LBS(V_0, \theta), F, C\}$

Coarse Shape Mesh Generation

Shape Sculpting and Texture Generation in the Implicit Space

Marching Cube & Get the Animatable 3D Avatar

a) Static Avatar Generation

$t_{motion}$ = "running"

Candidate Poses Generation

Reference-based Animation with Motion Prior

Final Animated 3D Avatars

b) Motion Generation

1) arguing

2) running

3) praying

(i)  (ii)  (iii)  Ours

(i)  (ii)  (iii)  Ours

(i)  (ii)  (iii)  Ours

a) Comparison with Baseline Methods

tired  sad  walking  squatting  raising both arms  washing hands  shooting basketball

b) More Results

(i) Direct optimization on SMPL parameter theta

(ii) Direct optimization on VPoser latent code

(iii) Multi-Modal RealNVP
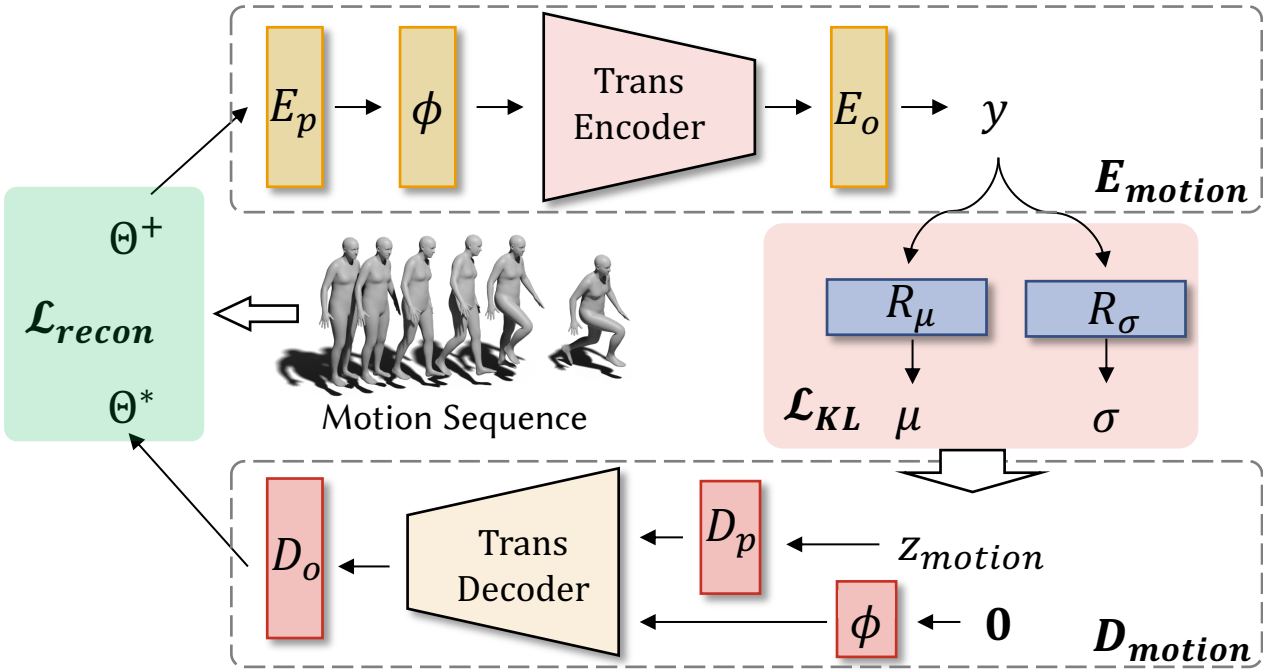
## A) MOTION VAE TRAINING

## B) CLIP-GUIDED OPTIMIZATION ON THE MOTION VAE



a) Reconstruction Loss $\mathcal{L}_{pose}$

b) Delta Loss $\mathcal{L}_{delta}$

c) CLIP Loss $\mathcal{L}_{clip}$

$$\mathcal{L}_{pose} = \sum_{i=1}^{k} \lambda_{pose}(i) min_j\{\|\theta_i - \Theta_j\|\}$$

$$\mathcal{L}_{delta} = -\sum_{i=1}^{L-1} \|\Theta_i - \Theta_{i+1}\|$$

$$s_i = 1 - \text{norm}(f_I) \cdot \text{norm}(f_T)$$

$$\mathcal{L}_{clip}^{m} = \sum_{i=1}^{L} \lambda_{clip}(i) \cdot s_i$$

$t_{motion}$ = "*raising both arms*"

# COMPARISONS OF MOTION GENERATION

Direct Interpolation          Direct motion VAE optimization (Baseline)          Ours

Brushing Teeth

# OVERALL RESULTS

An Overweight Man; Financial Manager; Excited
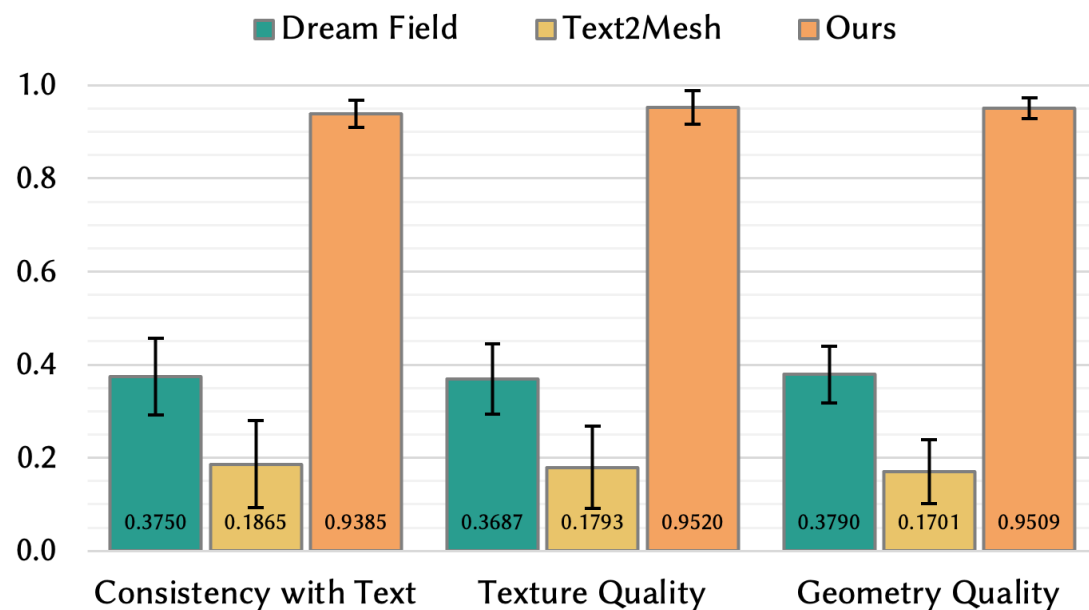
An Overweight Man; Sumo Wrestler; Sitting

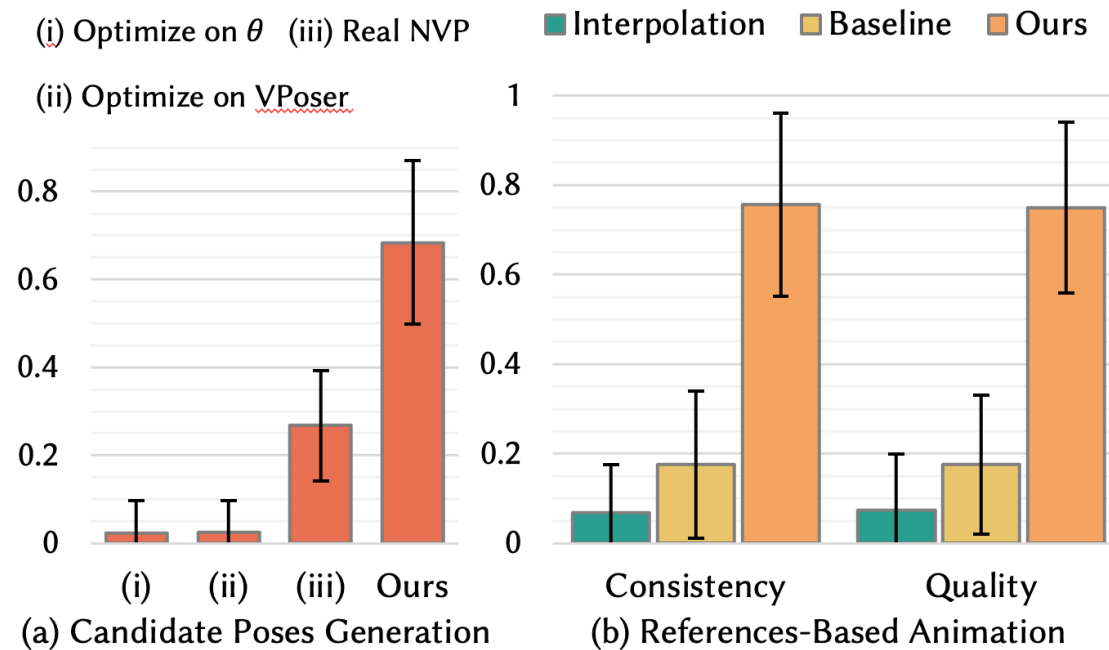A Strong Man; Firefighter; Kicking Soccer

## A) STATIC AVATAR GENERATION



## B) MOTION GENERATION

## LIMITATIONS

- Low quality of generate avatar.

- Small variations across different runs.

- Hard to generate out-of-distribution poses.

- Difficult to generate stylized motions.

## POTENTIAL NEGATIVE IMPACT

- Gender bias.

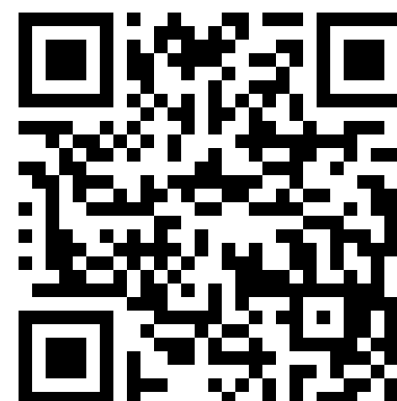- Misused to make fake videos of celebrities.



seed = 0    seed = 29    seed = 210    seed = 9176    seed = 12789        seed = 0    seed = 29    seed = 210    seed = 9176    seed = 12789

Teacher                                                        Student

(a) Text2Mesh

(b) Ours

# THANK YOU

## CODES ARE AVAILABLE

GitHub          Project Page