

AI-Driven Visual Content Generation

Ziwei Liu

刘子纬

Nanyang Technological University



NANYANG
TECHNOLOGICAL
UNIVERSITY
SINGAPORE

S-LAB
FOR ADVANCED
INTELLIGENCE

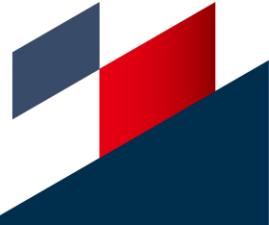
ABOUT

MMLab@NTU

MMLab@NTU was formed on the 1 August 2018, with a research focus on computer vision and deep learning. Its sister lab is [MMLab@CUHK](#). It is now a group with three faculty members and more than 40 members including research fellows, research assistants, and PhD students.

Members in MMLab@NTU conduct research primarily in low-level vision, image and video understanding, creative content creation, 3D scene understanding and reconstruction. Have a look at the overview of [our research](#). All publications are listed [here](#).

We are always looking for motivated PhD students, postdocs, research assistants who have the same interests like us. Check out the [careers](#) page and follow us on [Twitter](#).



AI-Generated Content (AIGC)



Movie



Game



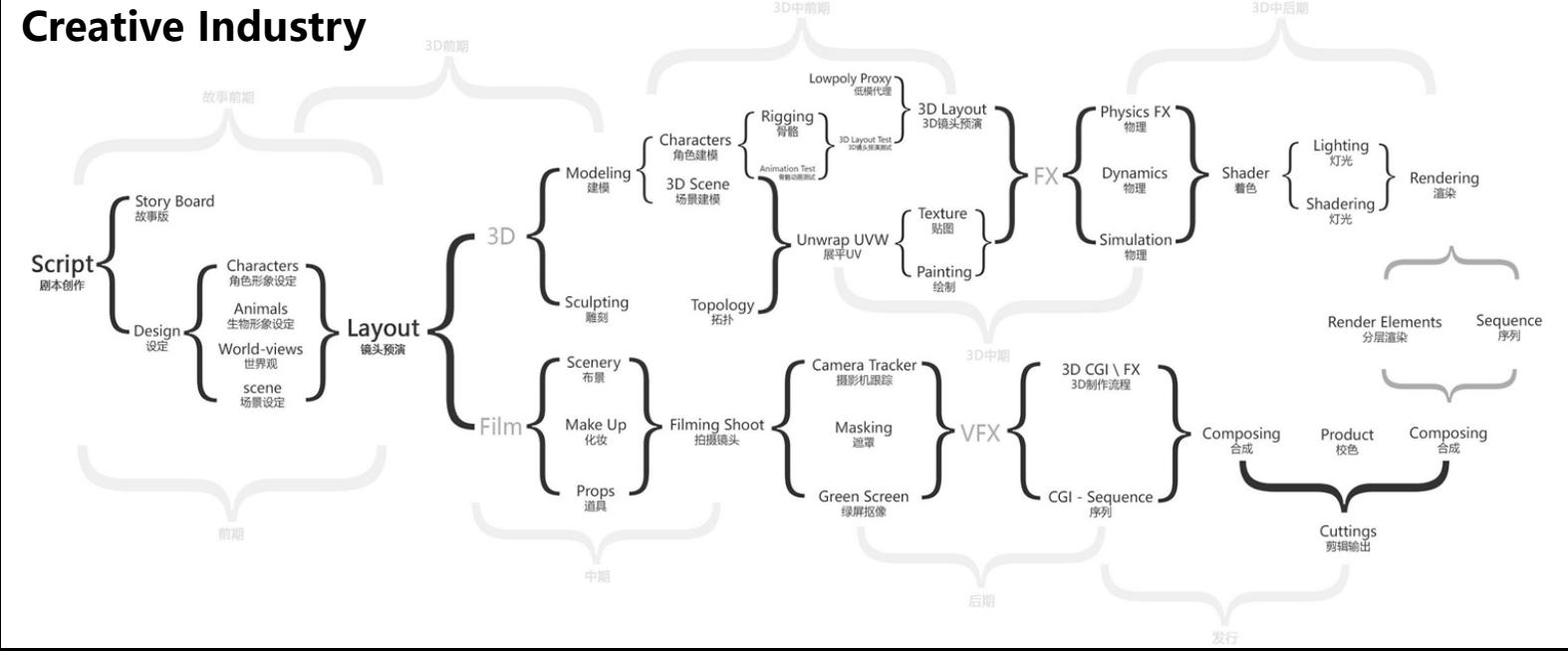
Anime



VTuber



Virtual Beings





2D Generation



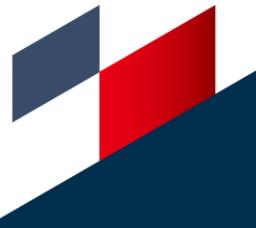
Motion Generation



3D Generation



Scene Generation

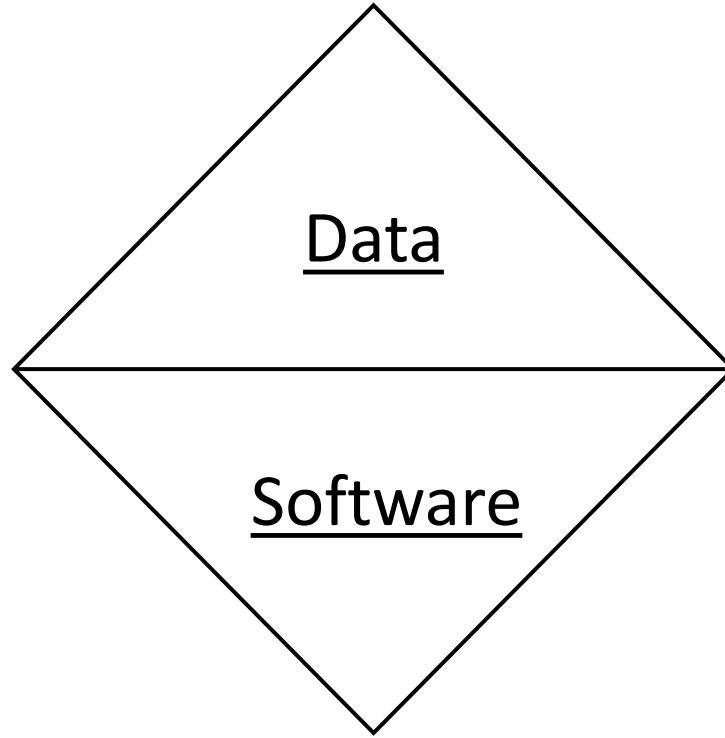




2D Generation



Motion Generation



3D Generation



Scene Generation

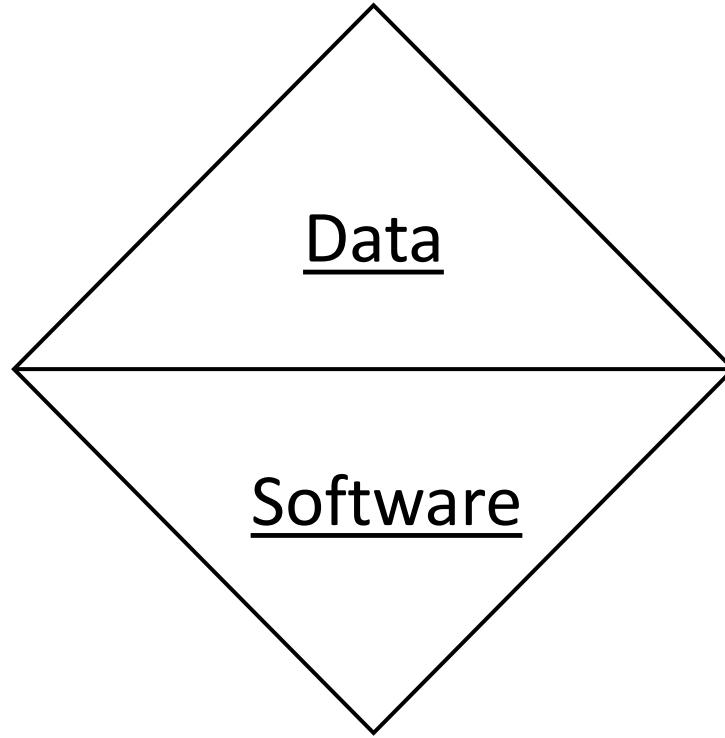




2D Generation



Motion Generation



3D Generation



Scene Generation

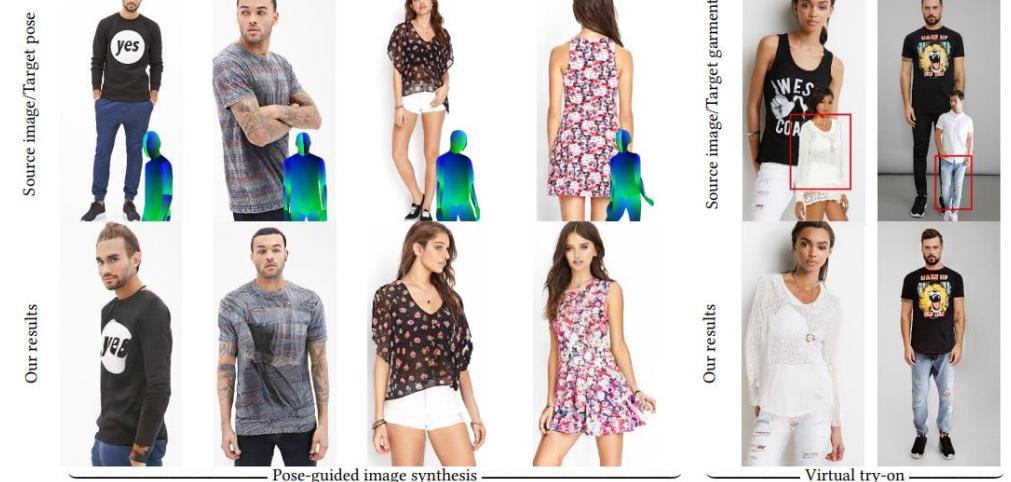


INTRODUCTION

- Human full-body images



- Pose Transfer
- Virtual try-on



Pose with Style [Albahar et al. 2021]

Text2Human

Load Pose

Generate Parsing

Save Image

Generate Human



Describe the shape.



Waiting for the generated result.

Describe the textures.



Parsing Palette

 top

leggings

 skin

ring

 outer

belt

 face

neckwear

 hair

socks

 dress

tie

 headwear

necklace

 pants

earstuds

 eyeglass

bag

 rompers

glove

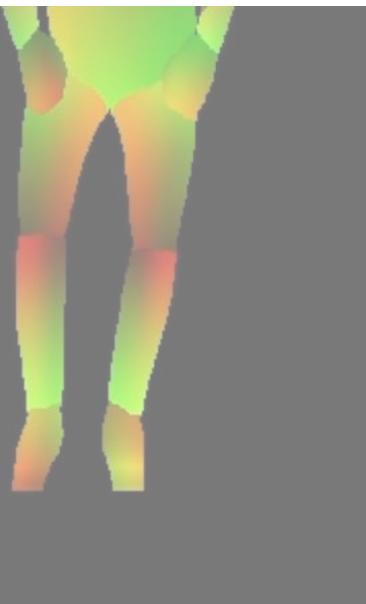
 footwear

background

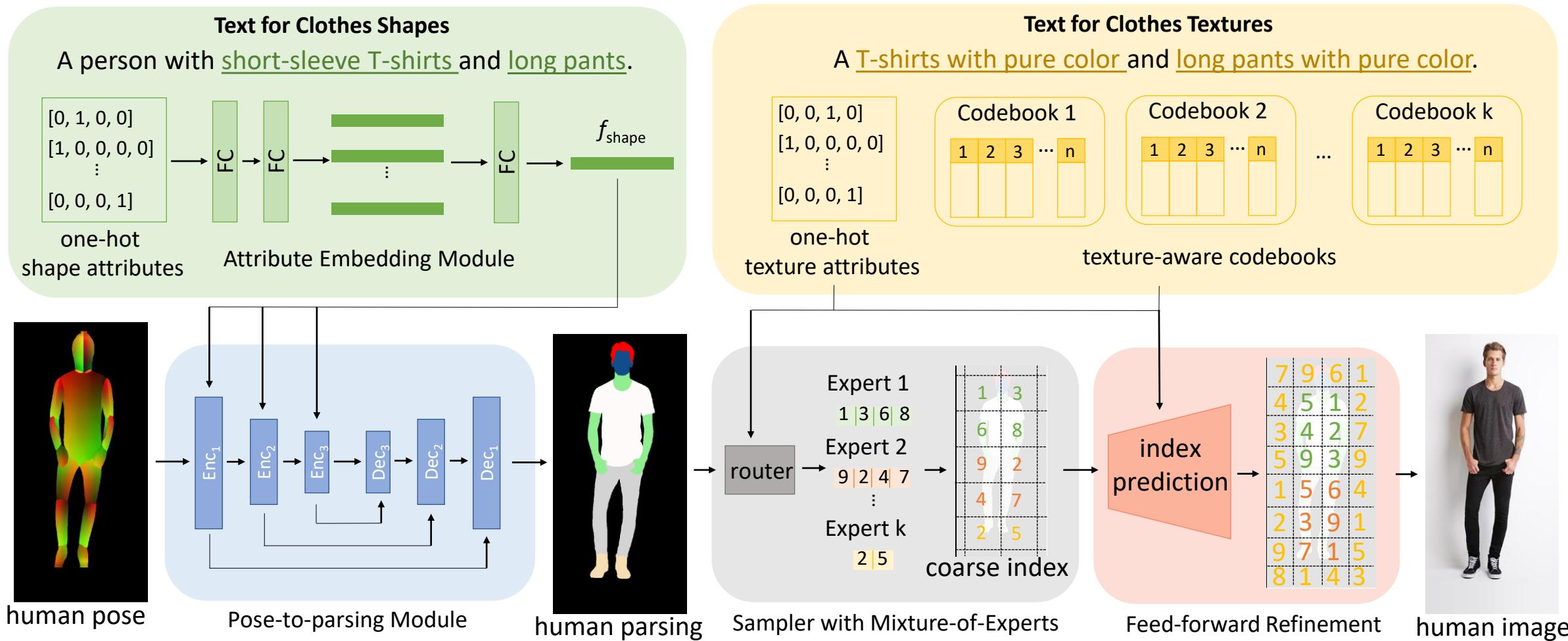
Provide the system with texts describing the shapes of clothes

Provide the system with texts describing the textures of clothes

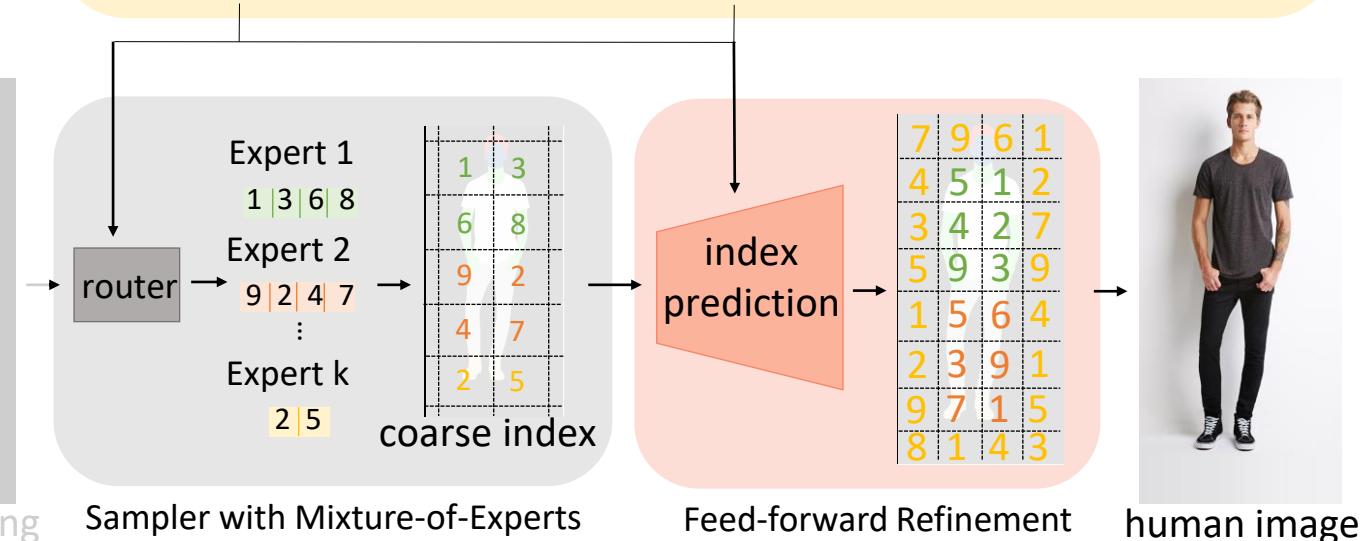
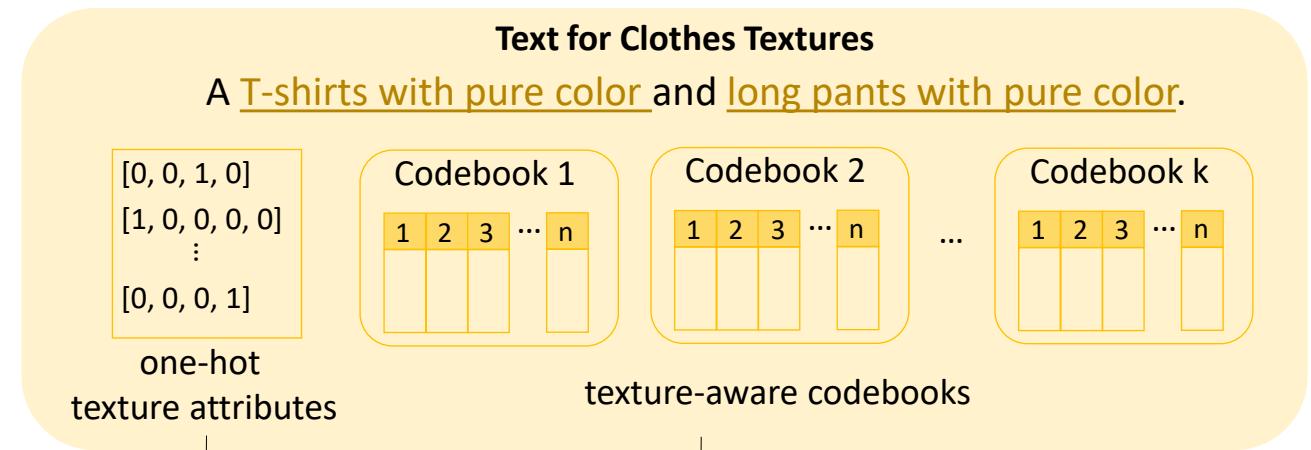
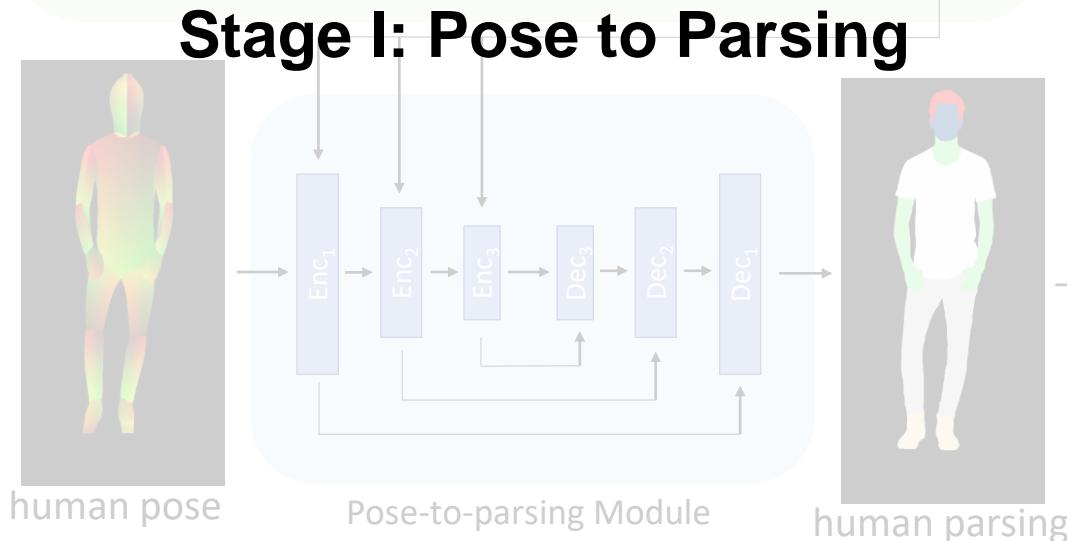
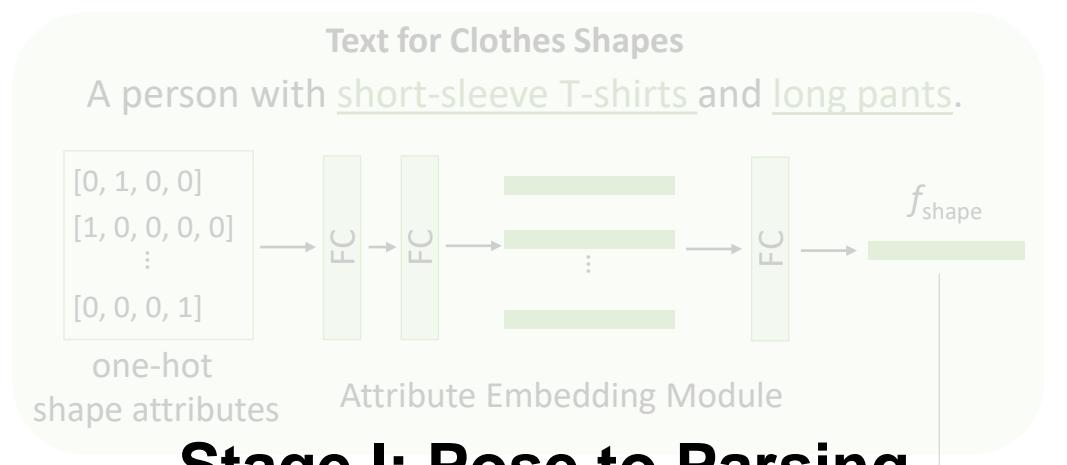
We propose a text-driven controllable human image generation task.



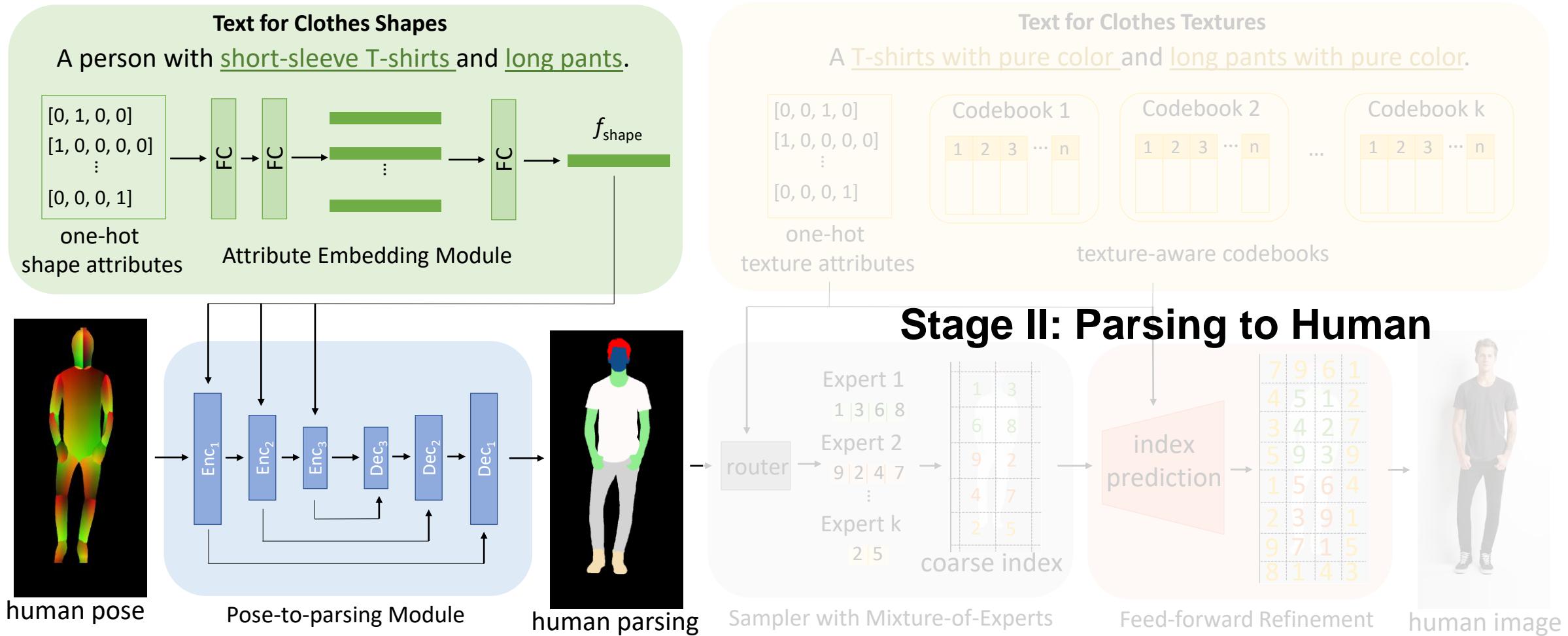
FRAMEWORK OF TEXT2HUMAN



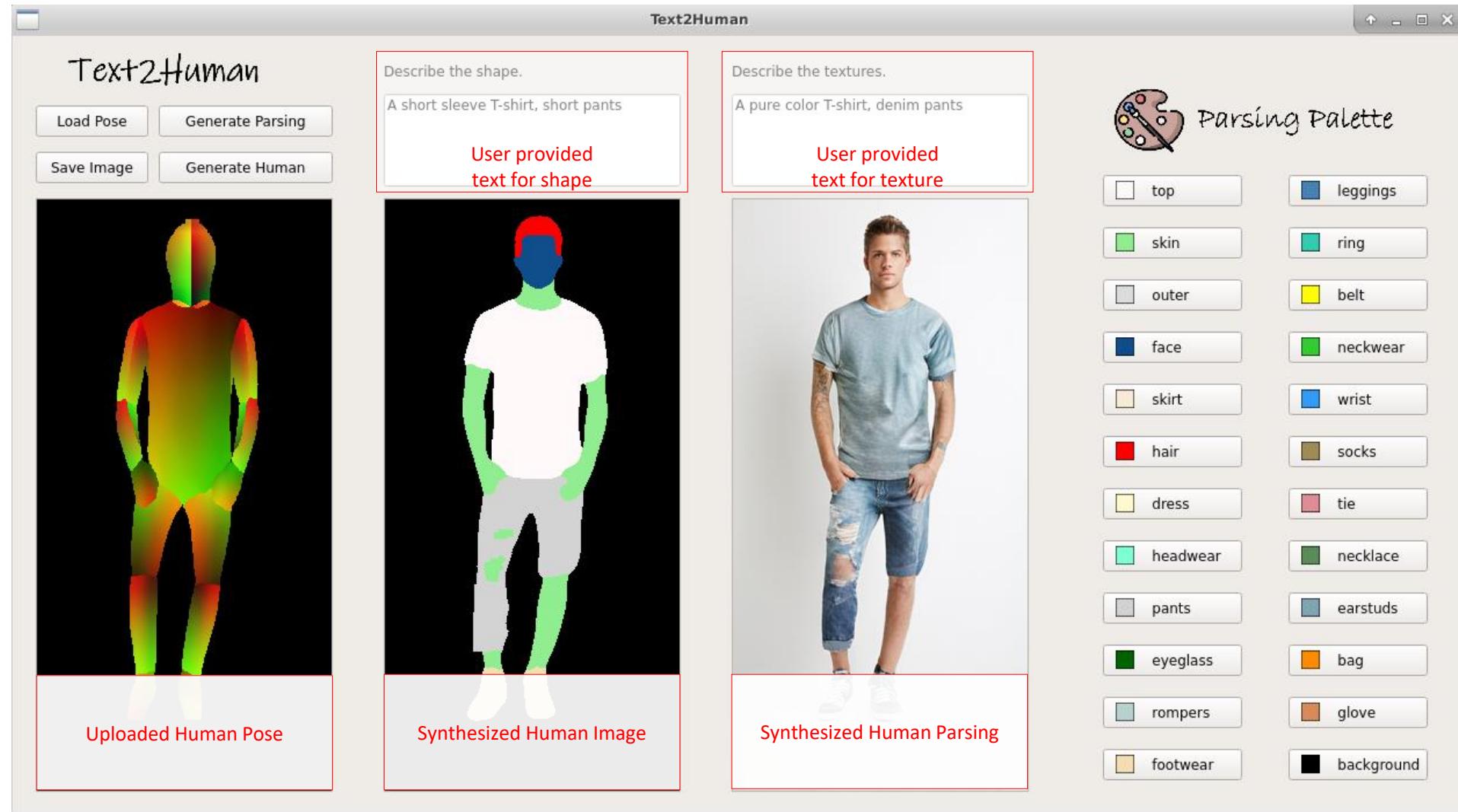
FRAMEWORK OF TEXT2HUMAN



FRAMEWORK OF TEXT2HUMAN



INTERACTIVE USER INTERFACE



DEEPFASHION-MULTIMODAL DATASET



DEEPFASHION-MULTIMODAL DATASET

- 44,096 high-resolution human images, including 12,701 full body human images
- **manually annotated** the human parsing labels
- DensePose for each human image
- **manually annotated** the keypoints
- **manually annotated** with attributes
- textual description



human image



human parsing



densepose



key points

shapes:
sleeve length: sleeveless
lower length: three-point
...
hat: no
socks: no
wrist accessory: yes
belt: no
neckline: suspenders
neckwear: no

Textures:
upper: cotton, graphic
lower: cotton, graphic
outer: NA.

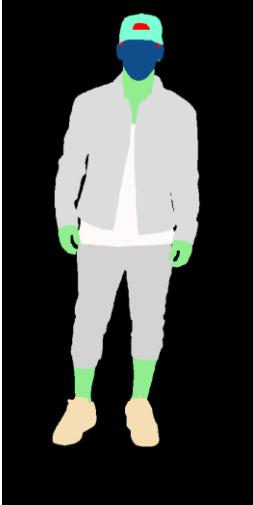
labels

The upper clothing has sleeves cut off, cotton fabric and graphic patterns. The neckline of it is suspenders. The lower clothing is of three-point length. The fabric is cotton, and it has graphic patterns. There is an accessory on her wrist.

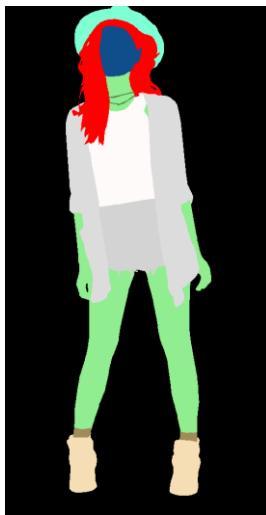
textual descriptions

EXPERIMENT

pure color upper clothes with a denim outer, seven-point and pure color pants



floral upper clothes with a pure color outer, three-point jeans



Parsing

Pix2PixHD

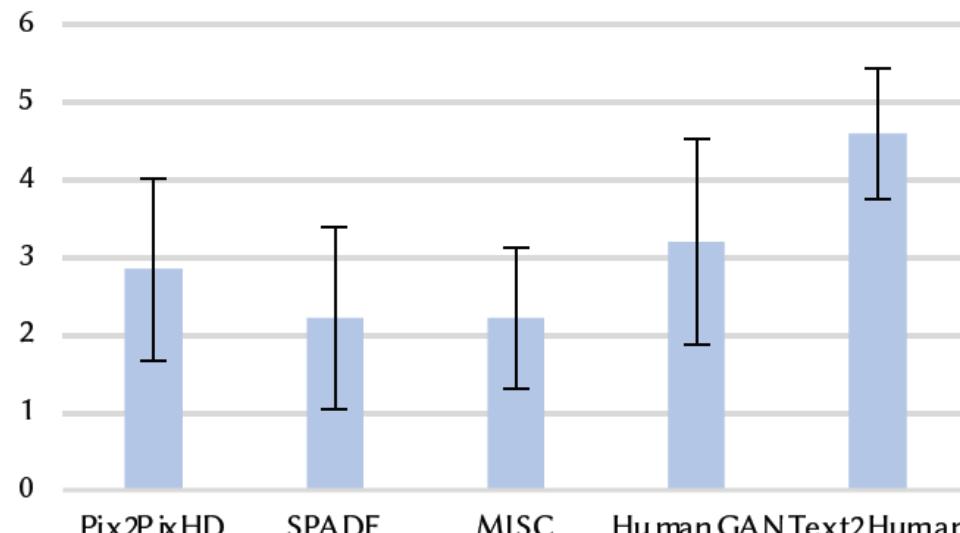
SPADE

MISC

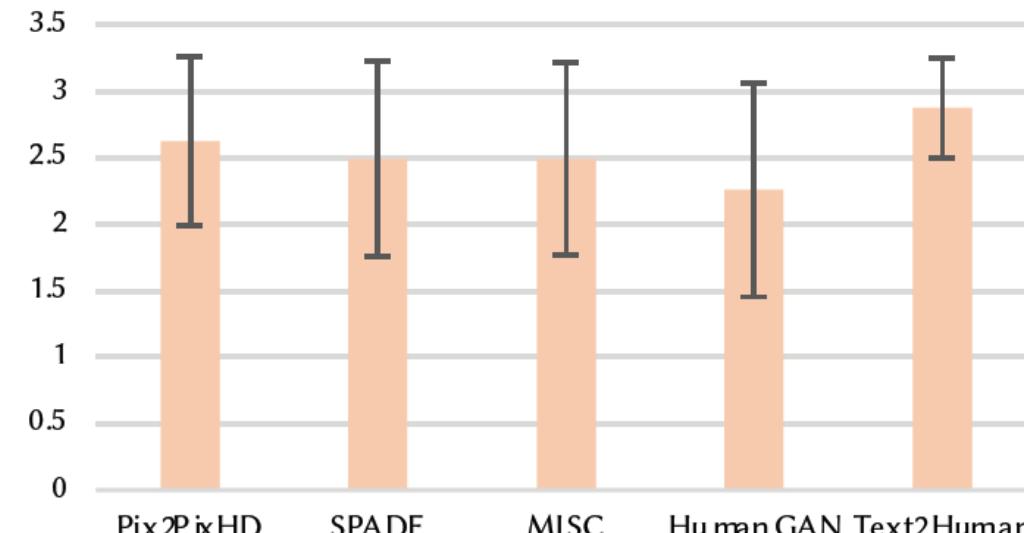
HumanGAN

Text2Human

EXPERIMENT



(a) photorealism



(b) texture consistency score

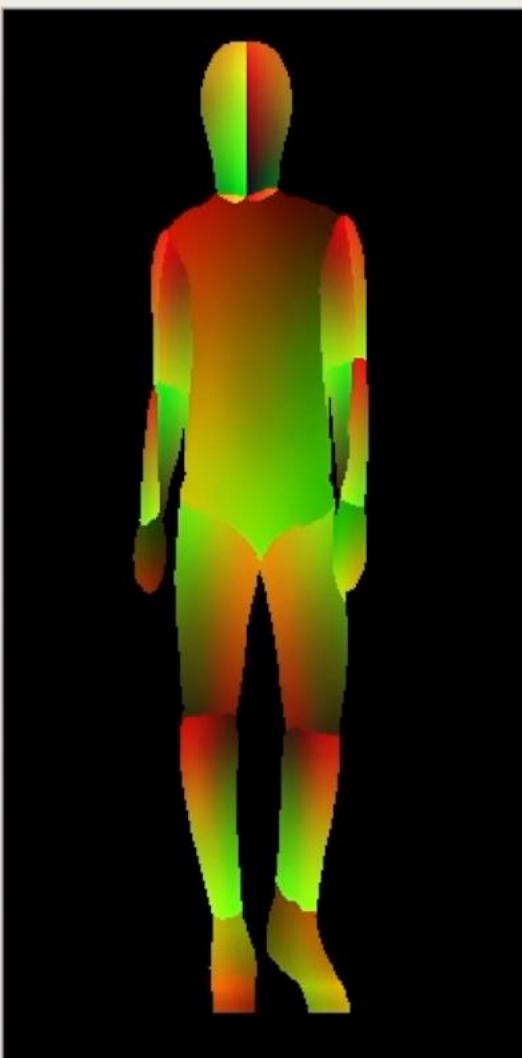
Text2Human

Load Pose

Generate Parsing

Save Image

Generate Human



Describe the shape.

 I

Describe the textures.



Parsing Palette

Modify the desired texture by texts.

 top leggings skin ring outer belt face neckwear skirt wrist hair socks dress tie headwear necklace pants earstuds eyeglass bag rompers glove footwear background

Text2Human

Describe the shape.

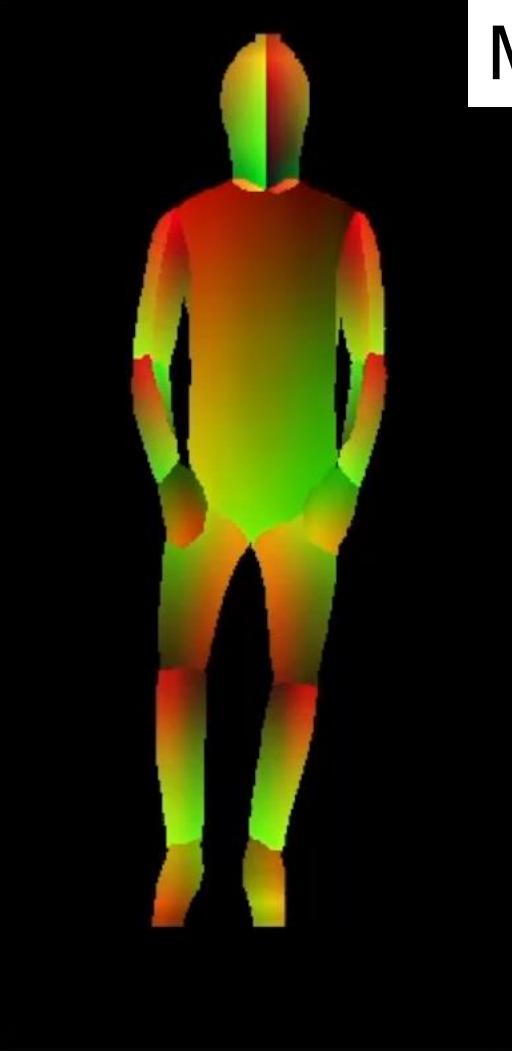
Describe the textures.

Load

Manually add more details on the generate human parsing clothes attributes

Save Image

Generate Human



Modify the desired clothes shape by texts.

<input type="checkbox"/> top	<input checked="" type="checkbox"/> leggings
<input checked="" type="checkbox"/> skin	<input type="checkbox"/> ring
<input type="checkbox"/> outer	<input type="checkbox"/> belt
<input checked="" type="checkbox"/> face	<input type="checkbox"/> neckwear
<input type="checkbox"/> skirt	<input type="checkbox"/> wrist
<input type="checkbox"/> hair	<input type="checkbox"/> socks
<input type="checkbox"/> dress	<input type="checkbox"/> tie
<input type="checkbox"/> headwear	<input type="checkbox"/> necklace
<input type="checkbox"/> pants	<input type="checkbox"/> earstuds
<input type="checkbox"/> eyeglass	<input type="checkbox"/> bag
<input type="checkbox"/> rompers	<input type="checkbox"/> glove
<input type="checkbox"/> footwear	<input type="checkbox"/> background

MORE SYNTHESIZED HUMAN IMAGES



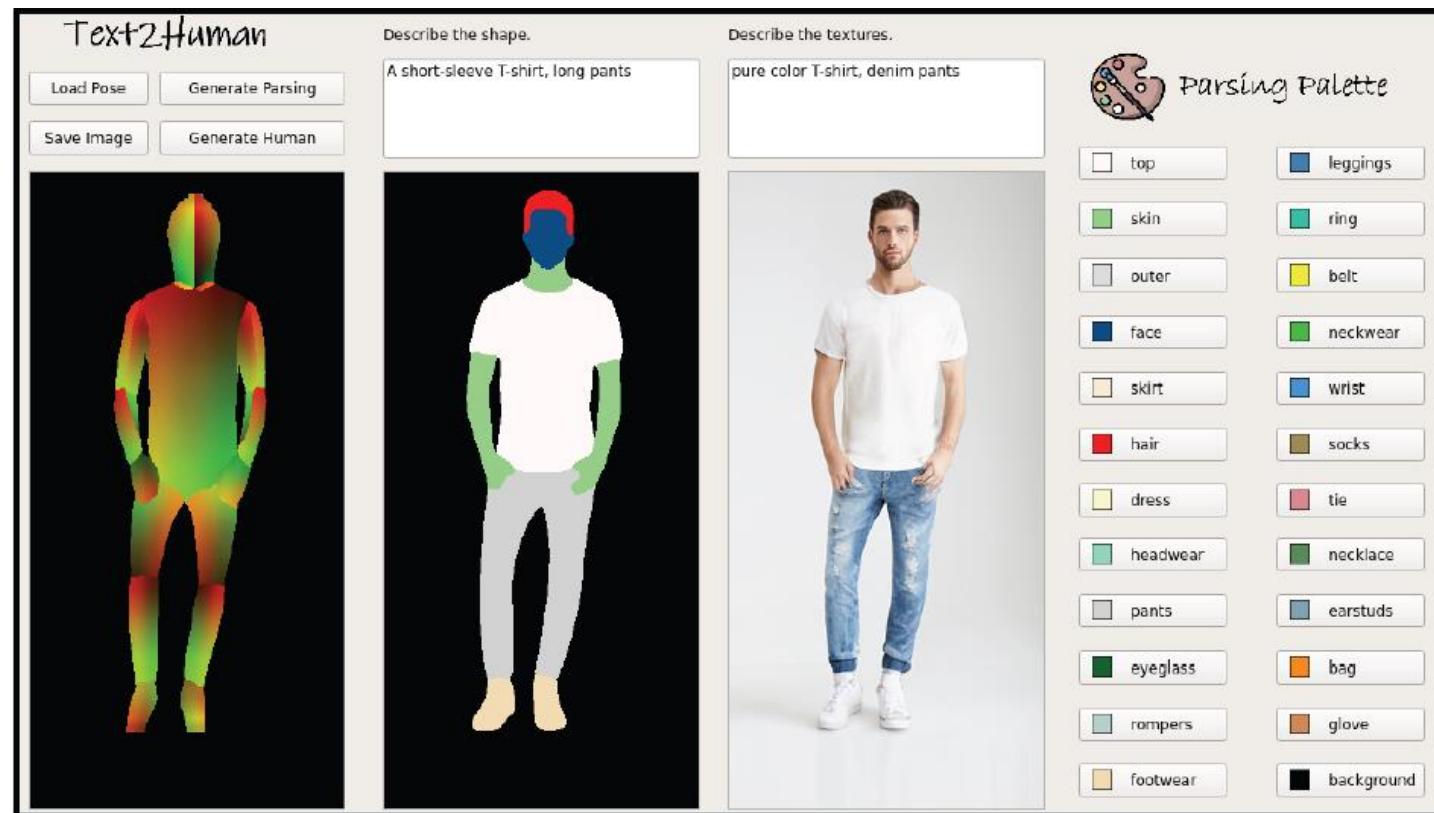






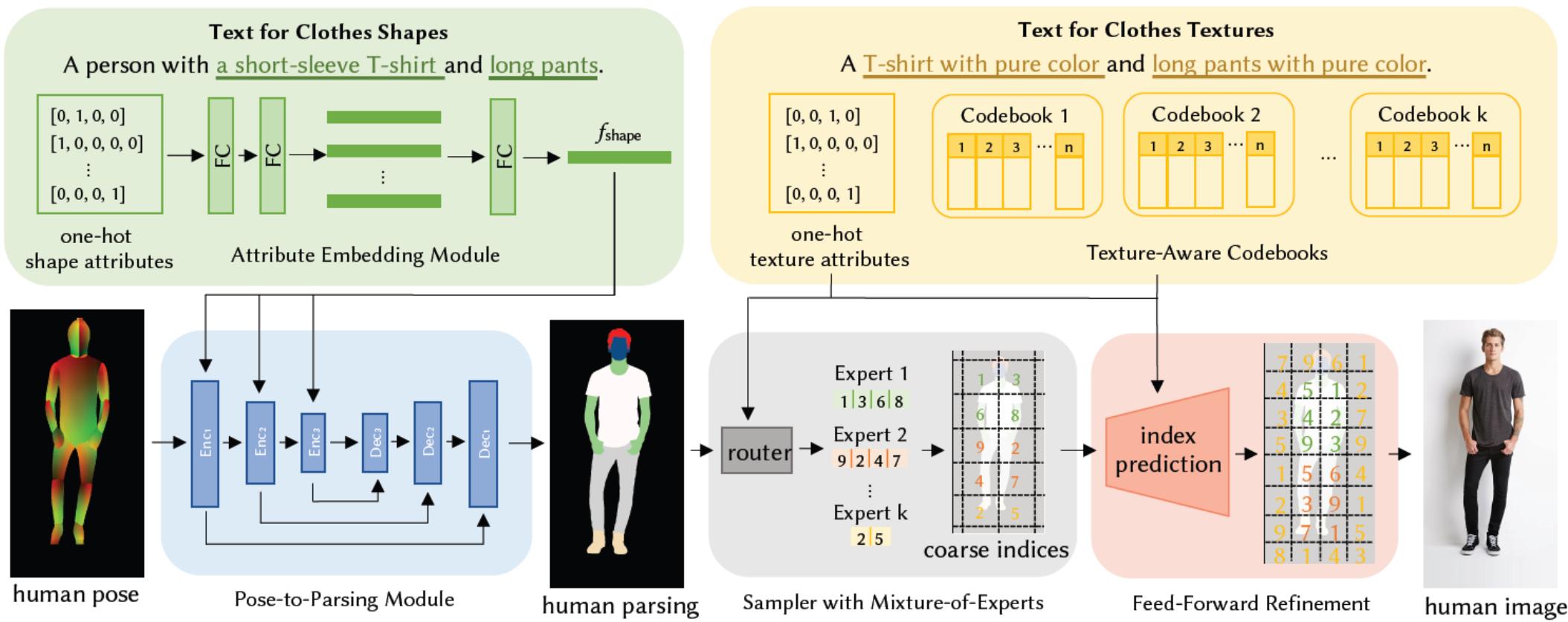
Task

Controllable Human Image Generation



Method

Text2Human



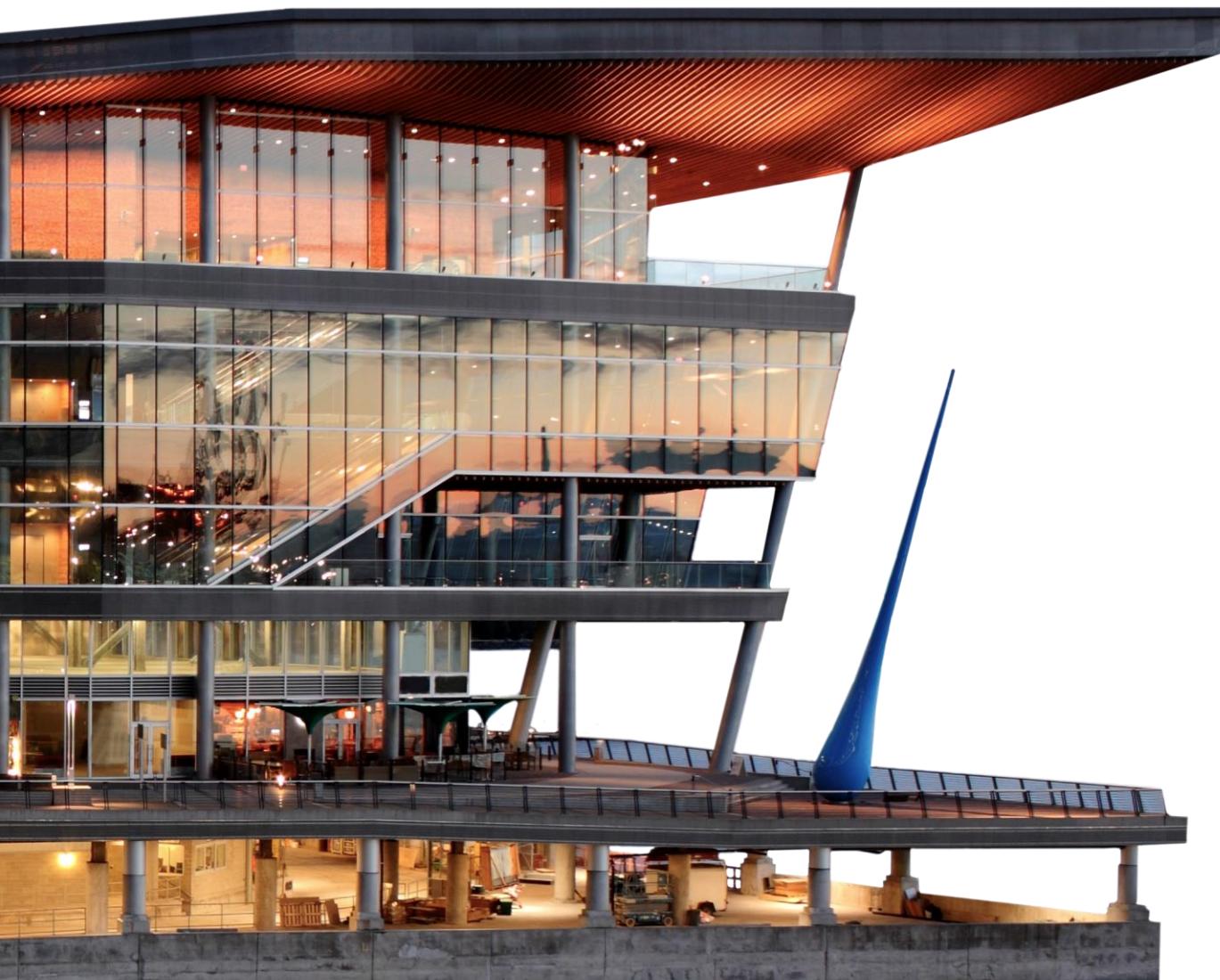
Dataset

DeepFashion-Multimodal





SIGGRAPH 2022
VANCOUVER+ 8-11 AUG



CODE AND MODELS

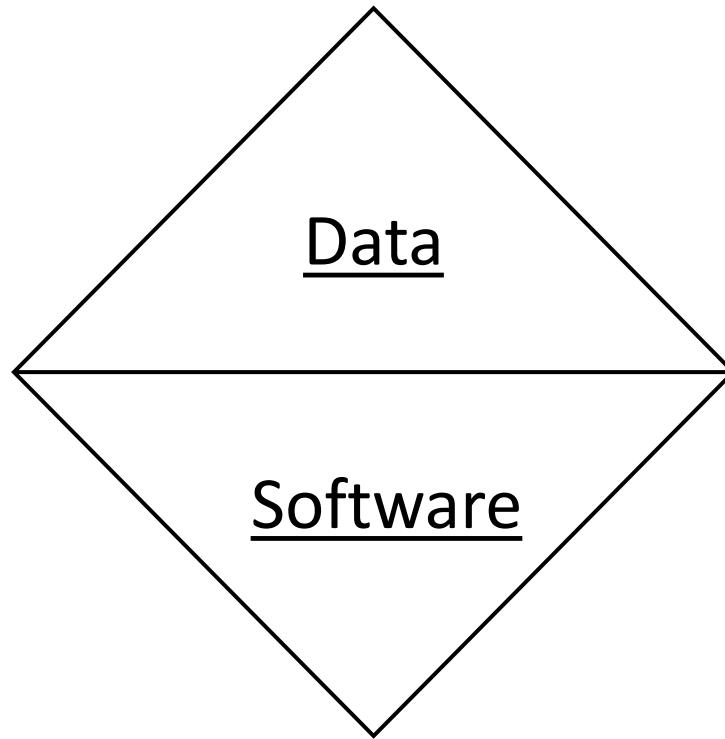




2D Generation



Motion Generation



3D Generation



Scene Generation



TEXT-DRIVEN IMAGE GENERATION



DALL·E [1]



DALL·E 2 [2]

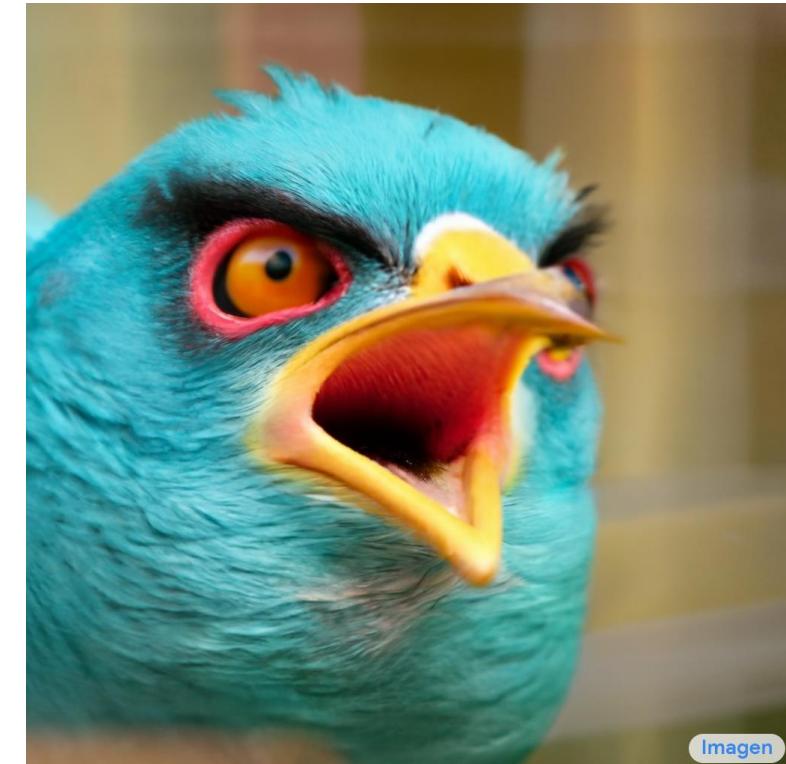


Imagen [3]

[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



Dream Field [1]

[1] <https://ajayj.com/dreamfields>

[2] <https://threedle.github.io/text2mesh/>



Text2Mesh [2]

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.



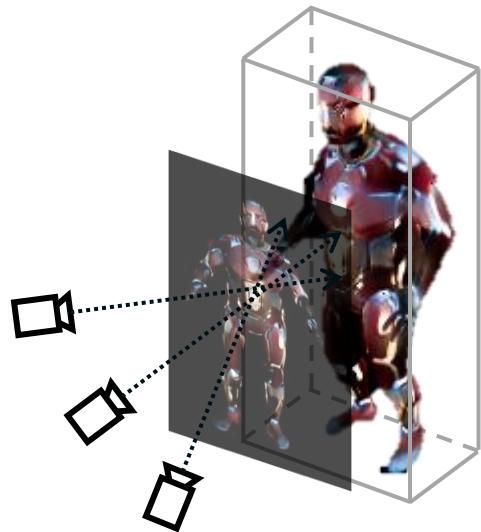
I want to generate a tall and skinny female soldier that is arguing.



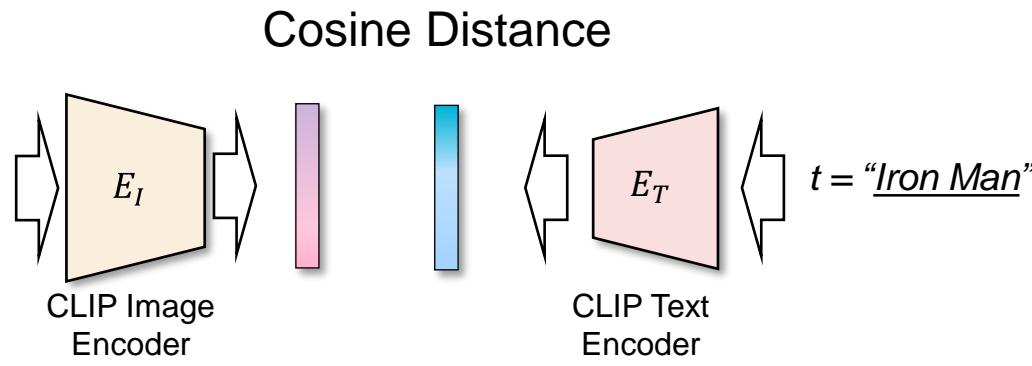
I want to generate an overweight sumo wrestler that is sitting.

TEXT-DRIVEN 3D GENERATION

CLIP + DIFFERENTIABLE RENDERING



a) Differentiable Rendering



b) Optimization guided by CLIP

AVATARCLIP: HOW IT WORKS

A) STATIC AVATAR GENERATION

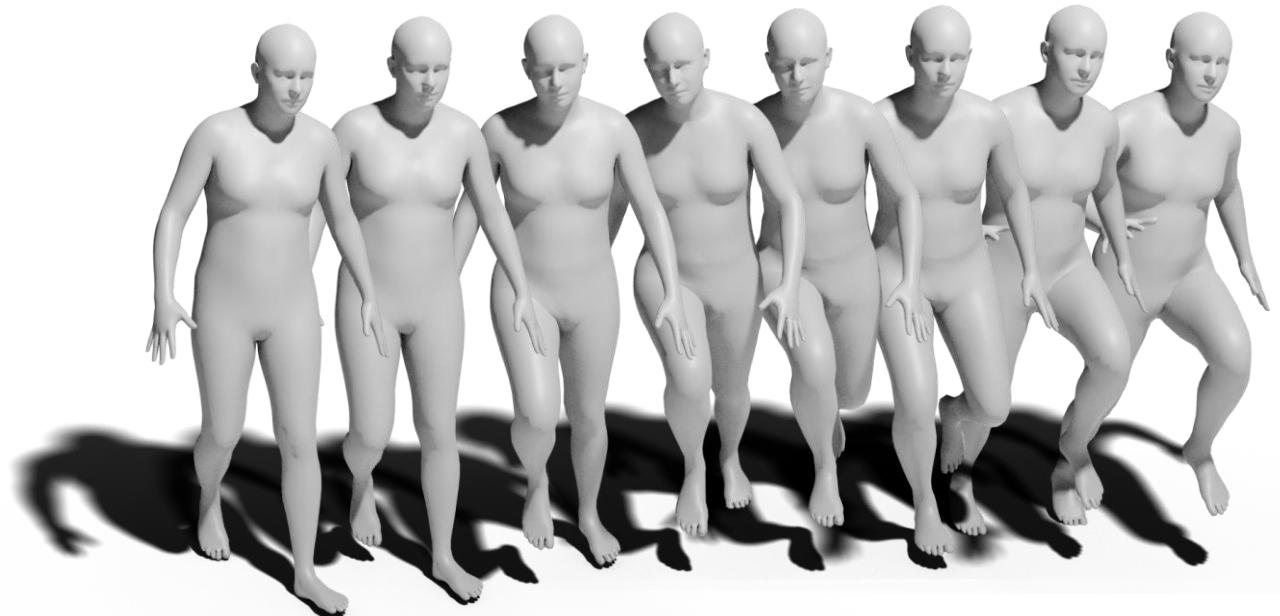
Shape Description: “*a tall and fat man*”

Appearance Description: “*Iron Man*”

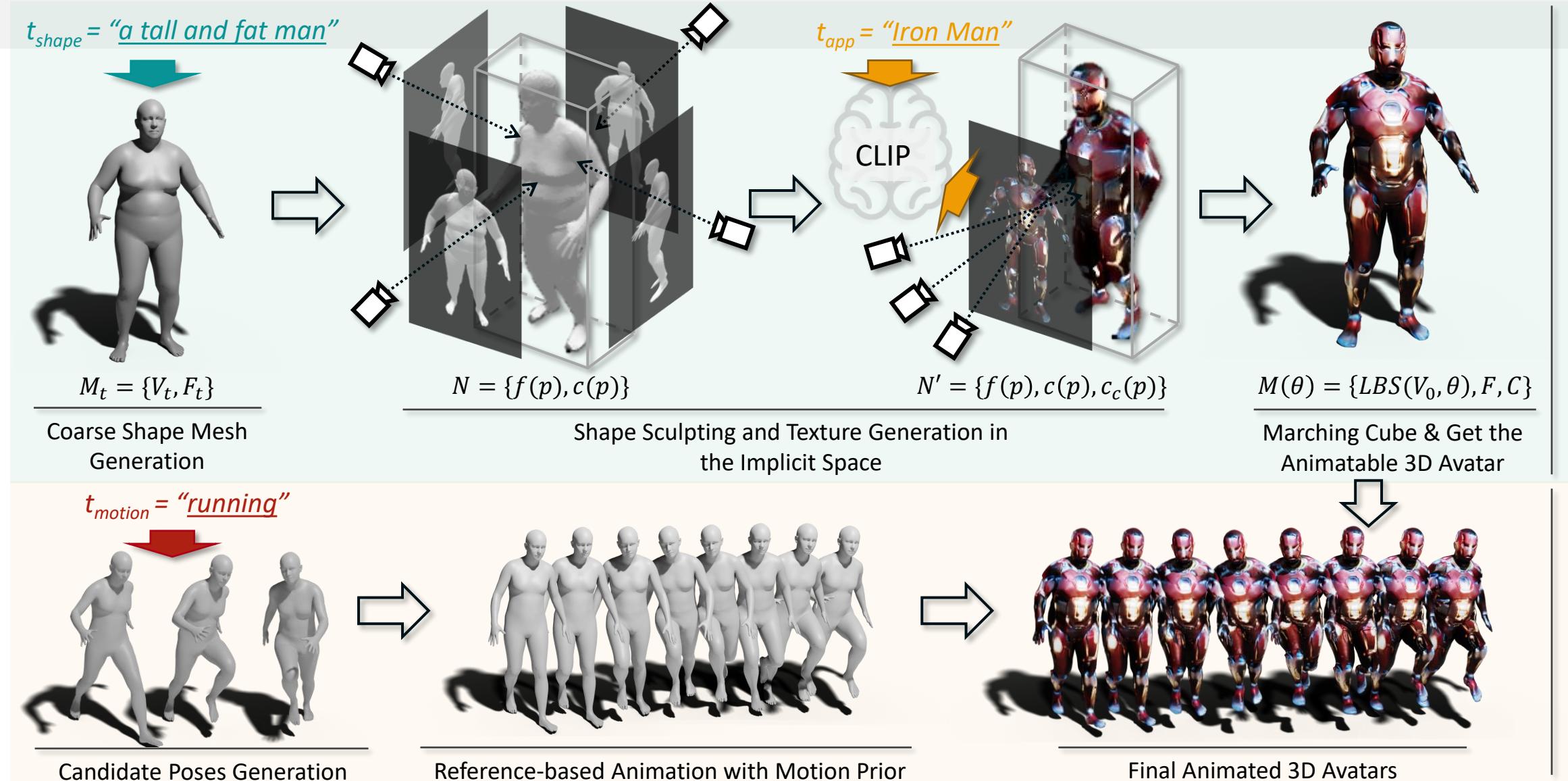


B) MOTION GENERATION

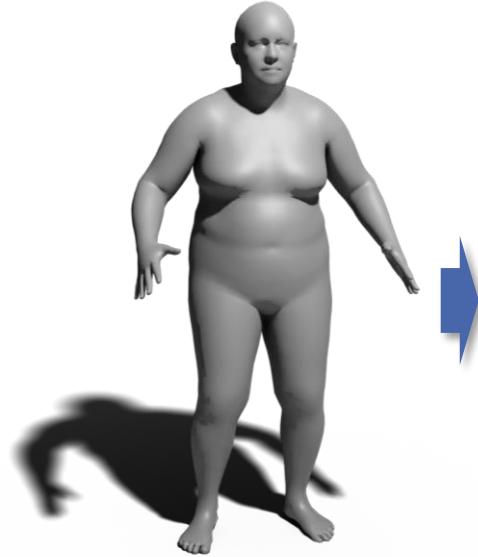
Motion Description: “*running*”



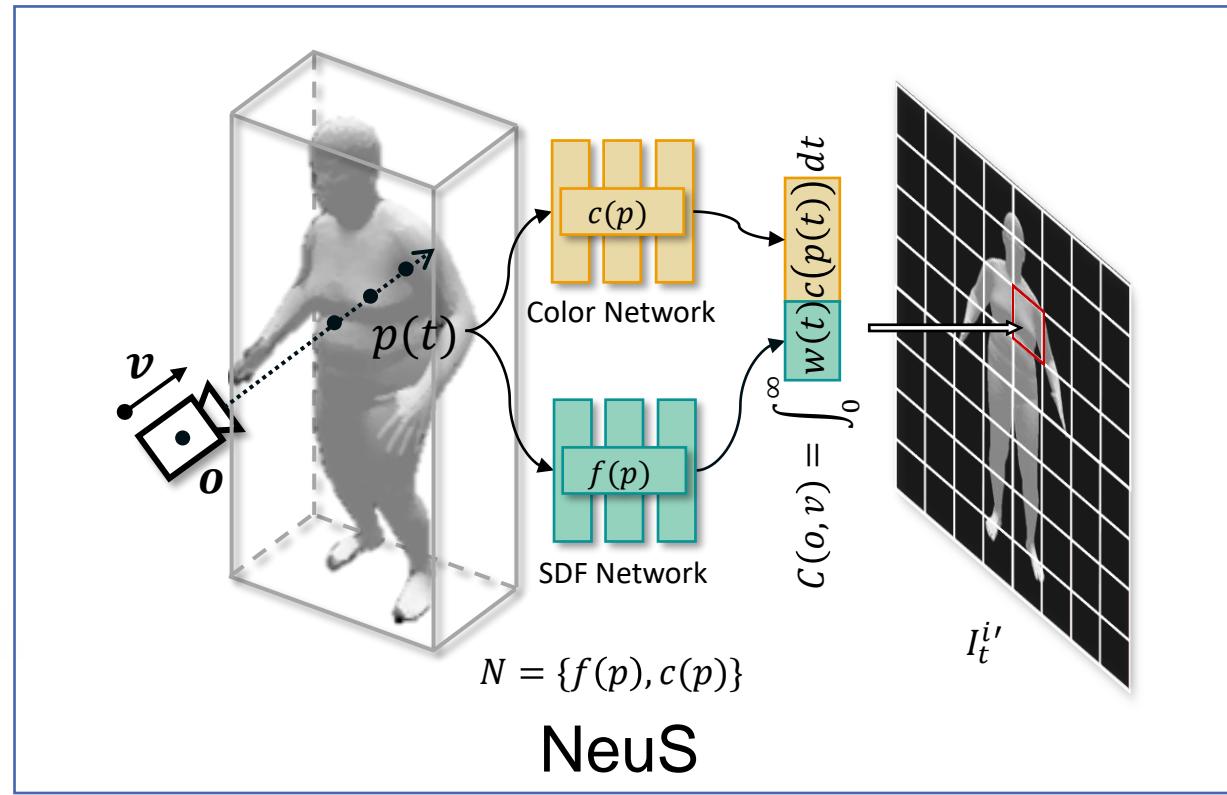
AVATARCLIP: DETAILED PIPELINE



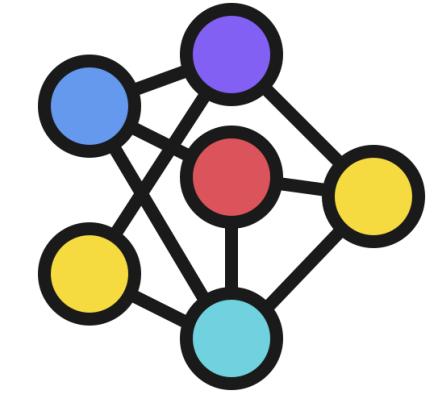
AVATARCLIP: TO THE IMPLICIT SPACE



Mesh

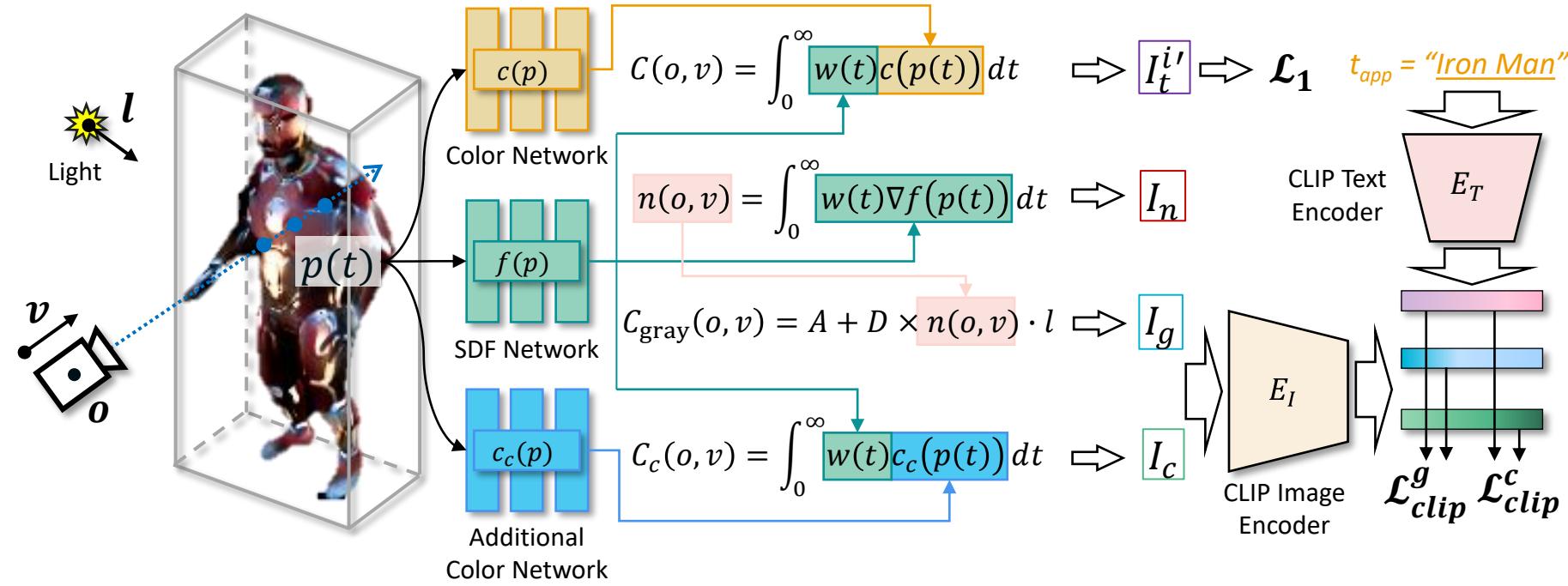


NeuS



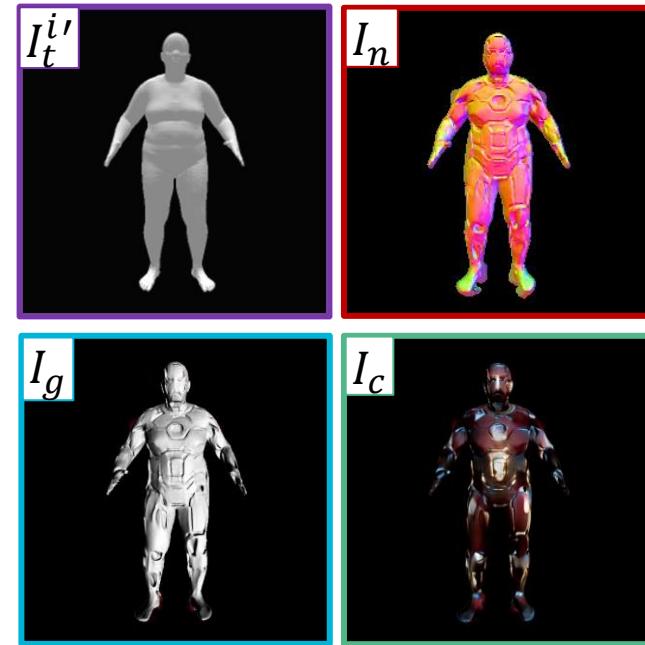
Implicit Function

AVATARCLIP: SHAPE SCULPTING AND TEXTURE GENERATION



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

b) Optimization



Examples of Intermediate Results

AVATARCLIP: OPTIMIZATION PROCESS

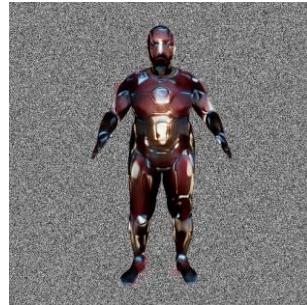
A) RANDOM BACKGROUND SEGMENTATION



1) Black



2) White



3) Gaussian
Noise



4) Chess
Board

B) SEMANTIC-AWARE PROMPT AUGMENTATION



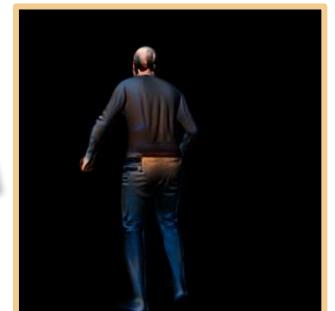
"the face of Steve Jobs"



"Steve Jobs"

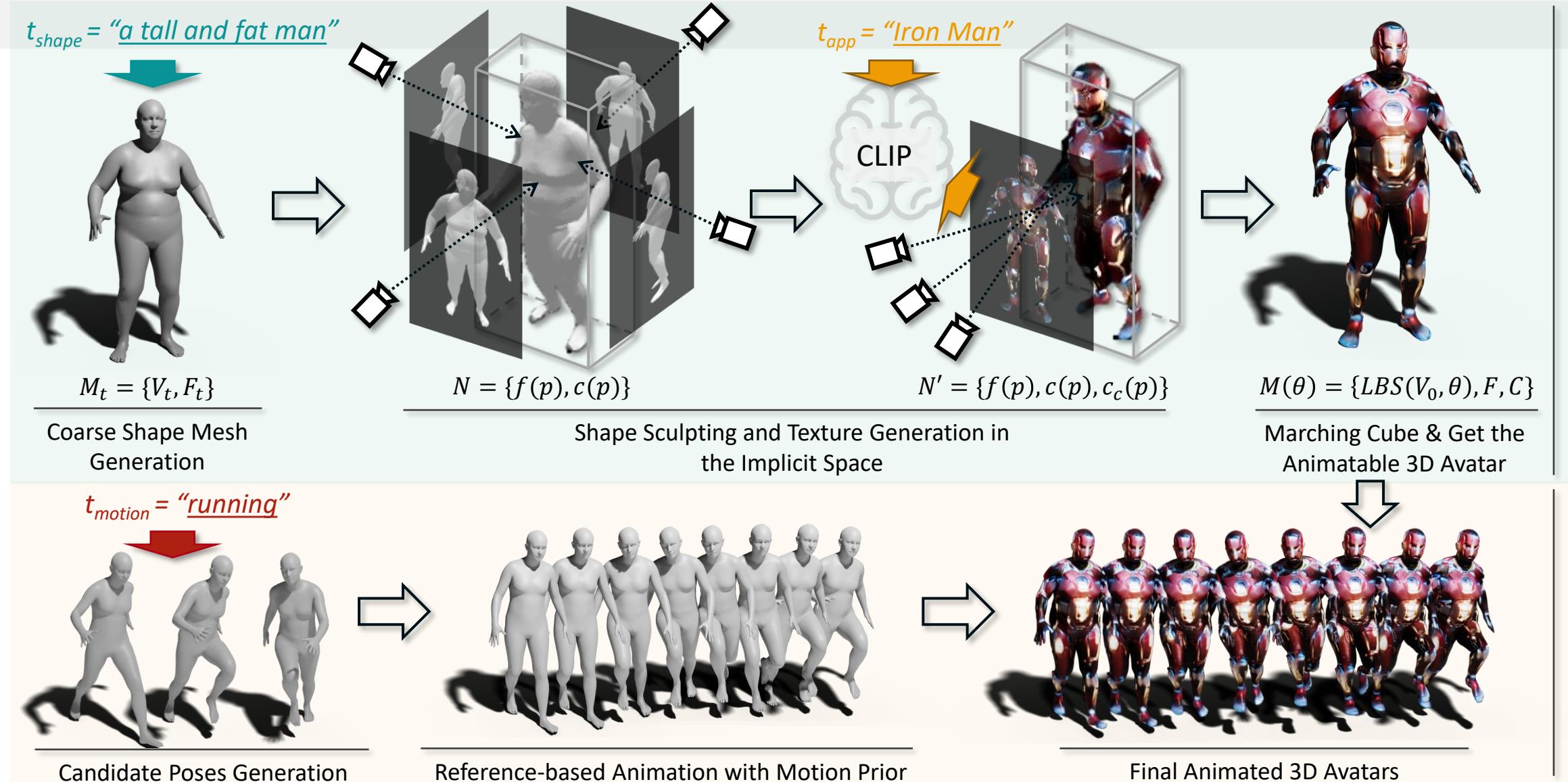


Implicit 3D Avatar $N' = \{f(p), c(p), c_c(p)\}$



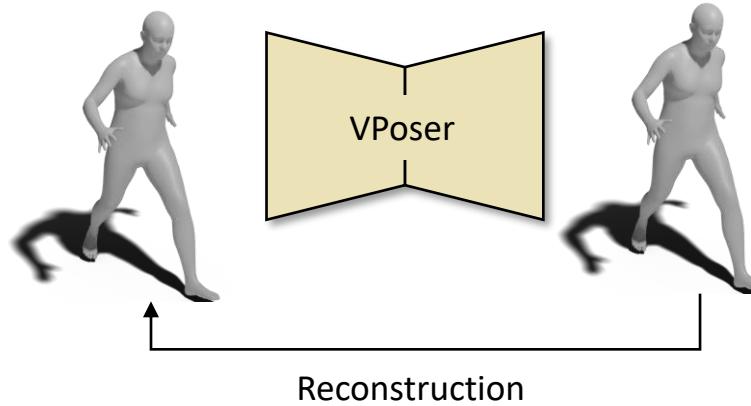
"the back of Steve Jobs"

AVATARCLIP: DETAILED PIPELINE

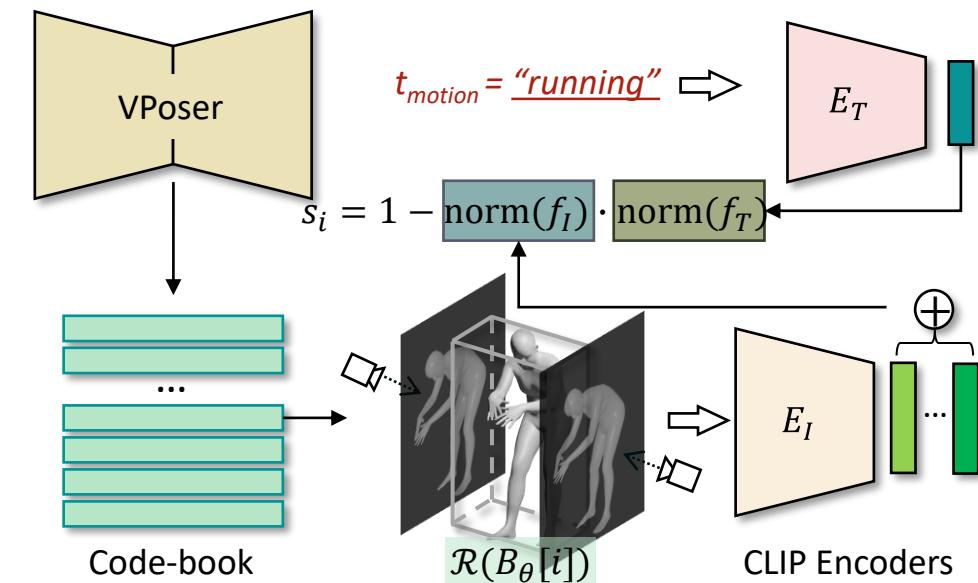


AVATARCLIP: CANDIDATE POSES GENERATION

A) POSE VAE (VPOSER)



B) CLIP-GUIDED CANDIDATE POSES QUERY



OVERALL RESULTS

AvatarCLIP

Create Your Own Avatar
with Natural Languages!



Describe the Shape



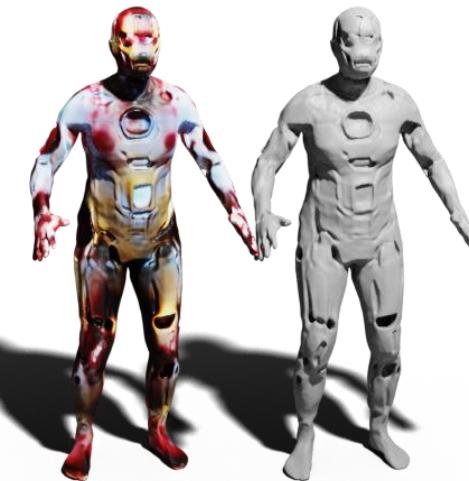
Generate

Next Step

Renderer Controller

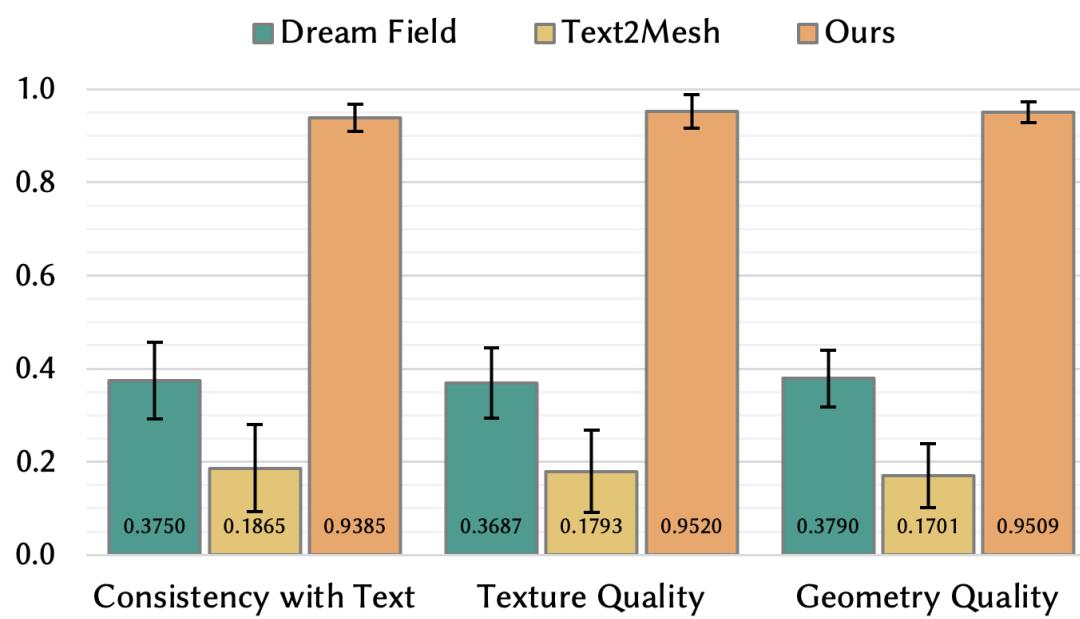
- Vertex Color
- Wireframe
- Normal

CONTROLLING & CONCEPT MIXING ABILITIES

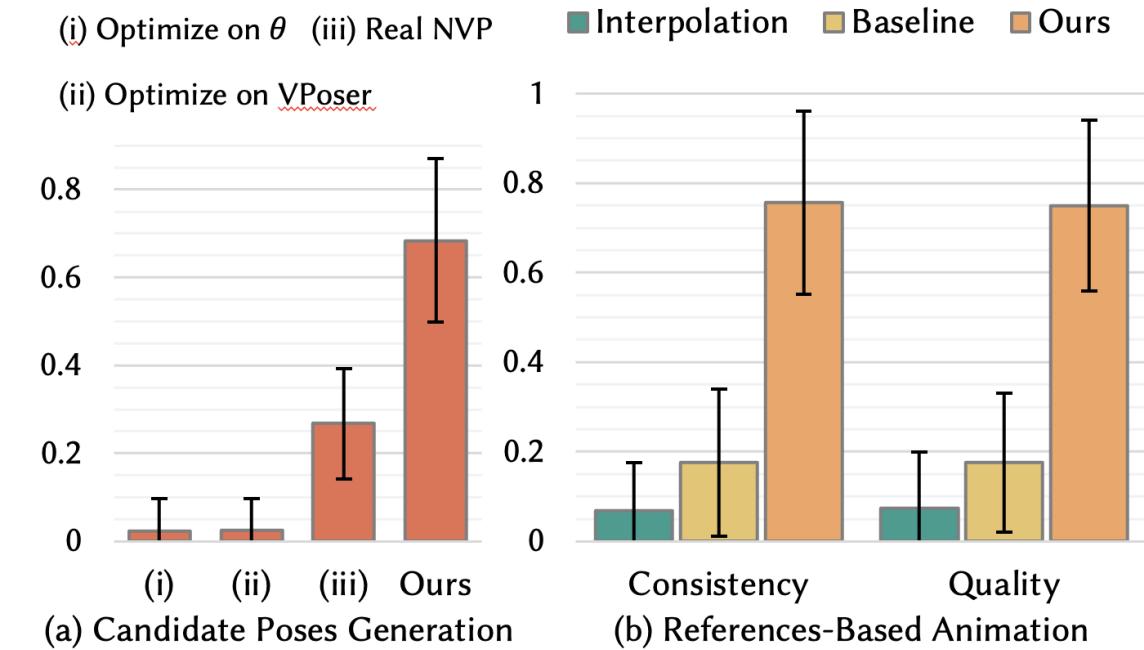


QUANTITATIVE RESULTS: USER STUDY

A) STATIC AVATAR GENERATION

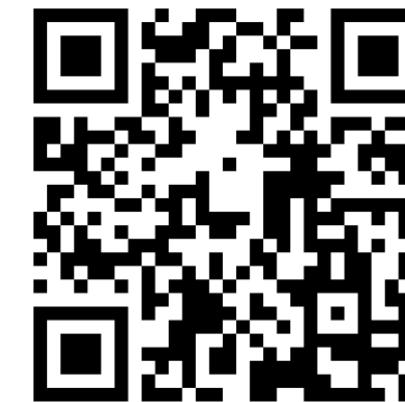


B) MOTION GENERATION





CODE AND MODELS



GitHub



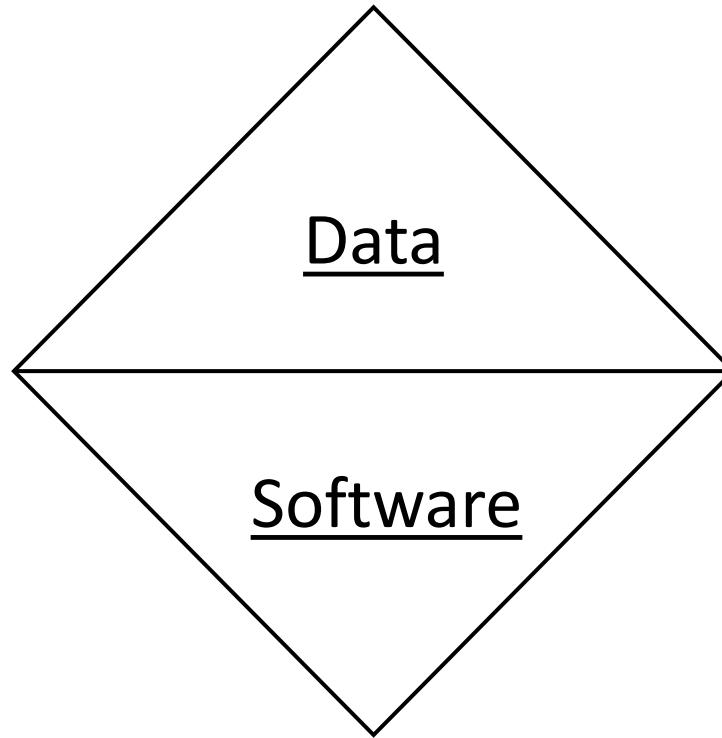
Project Page



2D Generation



Motion Generation



3D Generation



Scene Generation



3D Animation



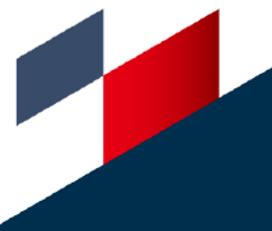
Video Games



Films



VTuber



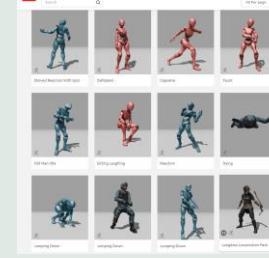
Motion Collection



Manual Editing

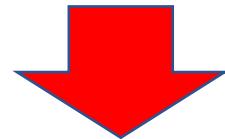


Motion Capture

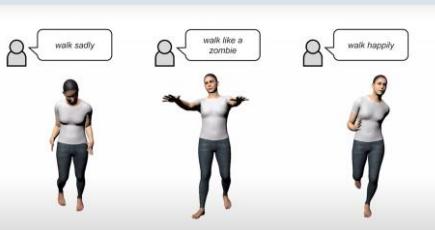


Gallery

1. **Expensive**
2. **Time-consuming**
3. **Not User-friendly**



Human Mesh Recovery

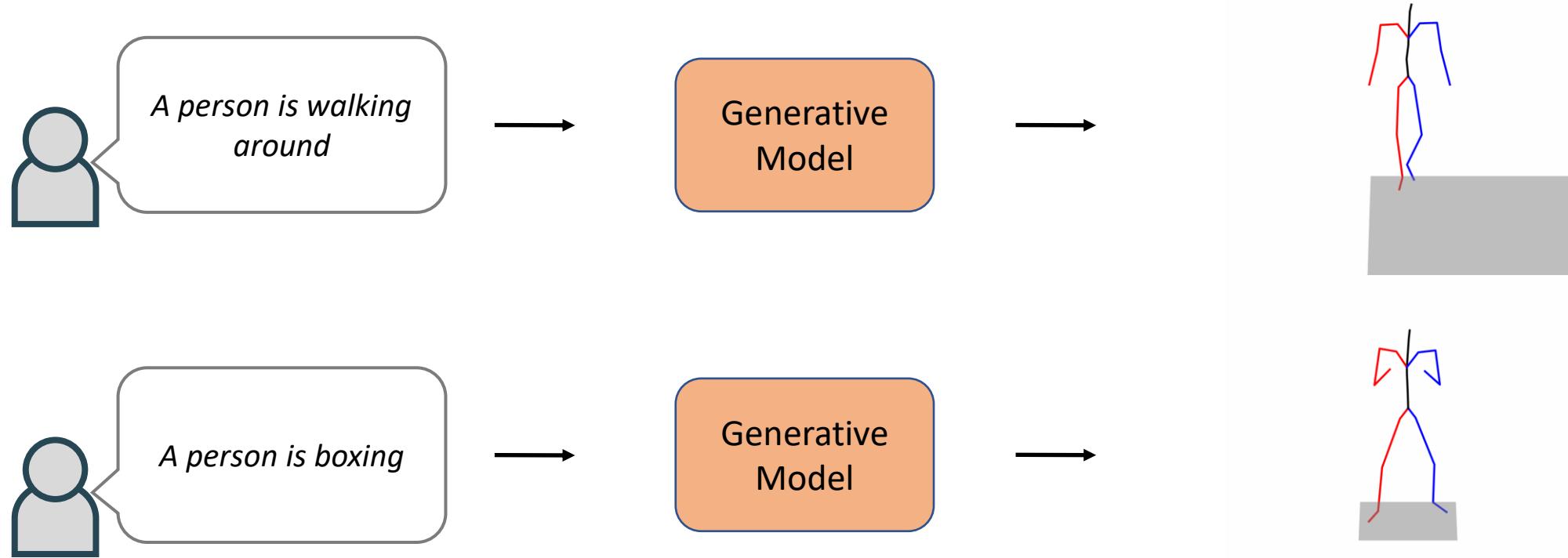


Conditional Motion Generation

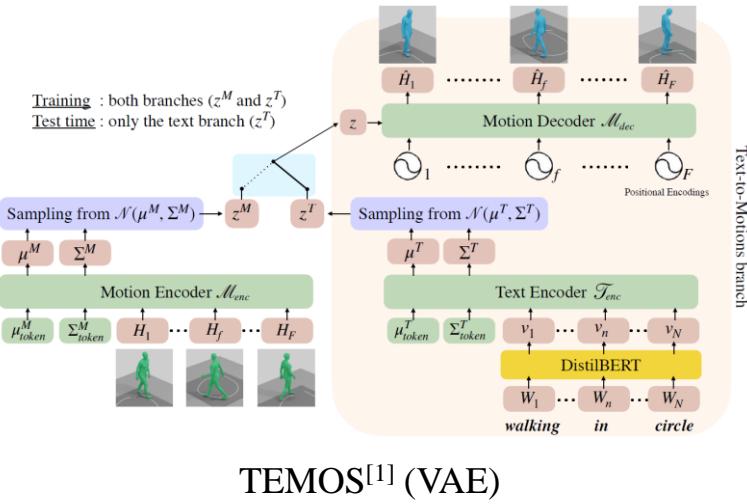
1. **Cheap**
2. **Efficient**
3. **User-friendly**



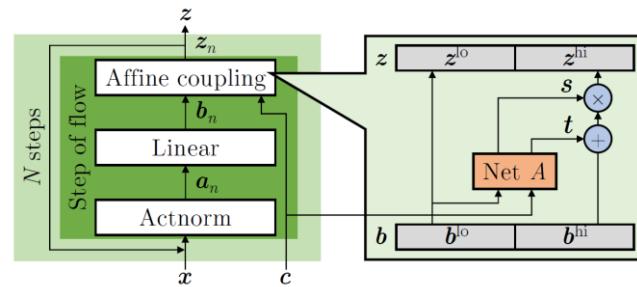
Text-driven Motion Generation



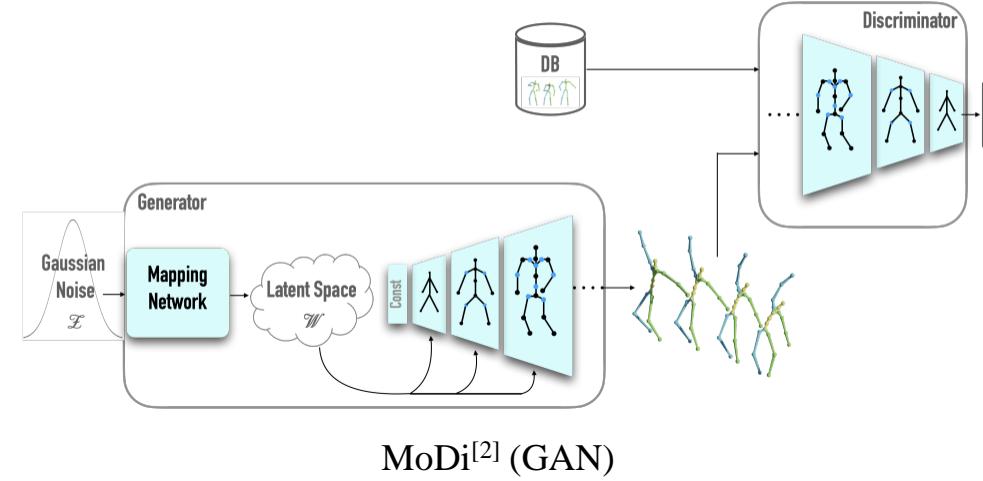
Classical Motion Generative Model



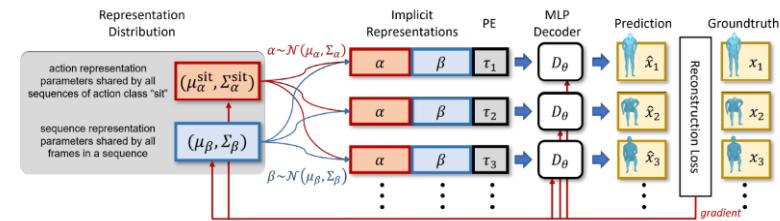
TEMOS^[1] (VAE)



MoGlow^[3] (Normalization Flow)



MoDi^[2] (GAN)



INR^[4] (Implicit Function)

Issues: 1) Hard to model complicated motion sequence 2) Lack of diversity

[1] Petrovich M, et al. Temos: Generating diverse human motions from textual descriptions. ECCV 2022

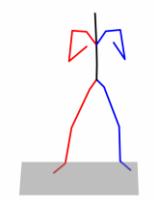
[2] Sigal R, et al. MoDi: Unconditional Motion Synthesis from Diverse Data. ArXiv 2022

[3] Henter GE, et al. Moglow: Probabilistic and controllable motion synthesis using normalising flows. TOG 2020

[4] Cervantes P, et al. Implicit neural representations for variable length human motion generation. ECCV 2022

Motion Generation with Diffusion Model

Diffusion Process



Add Noise



Add Noise

• • •

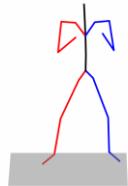
Add Noise



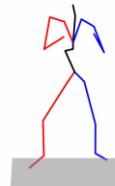
$$\mathbf{x}_0 \sim q(\mathbf{x}_0)$$

$$p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$$

Reverse Process



Denoise



Denoise

• • •

Denoise

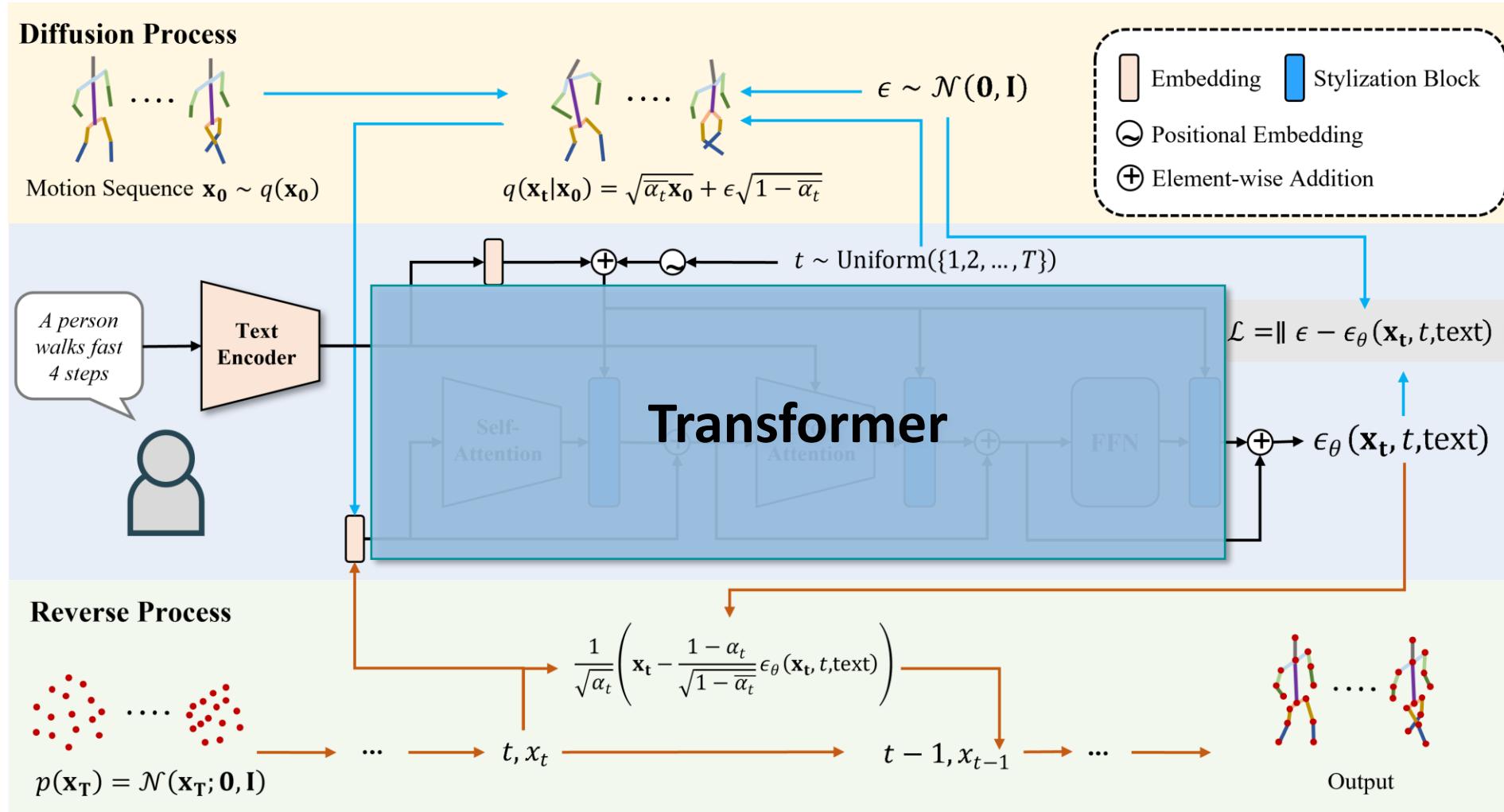


$$\mathbf{x}_0 \sim q(\mathbf{x}_0)$$

$$p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$$



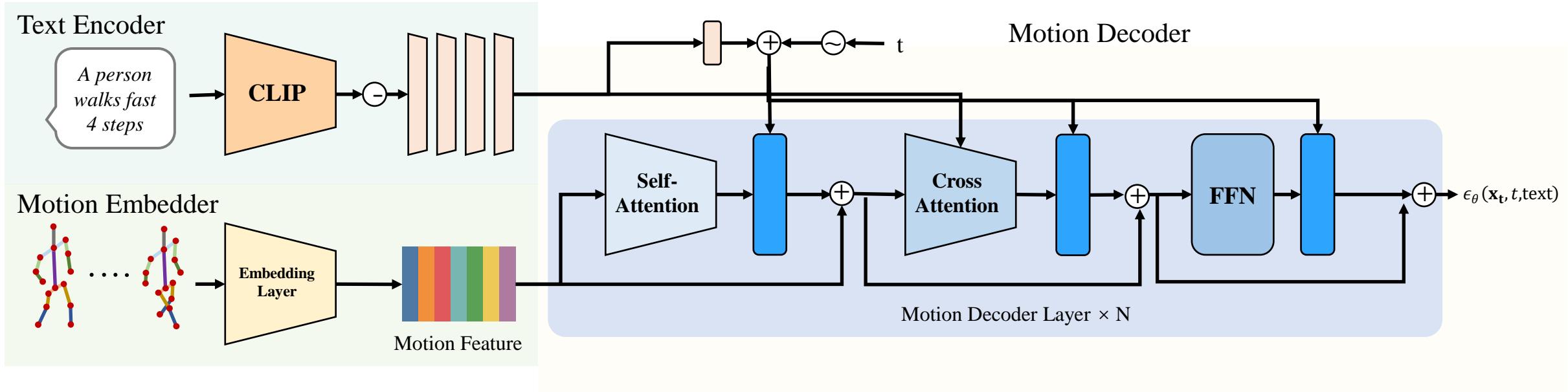
Framework



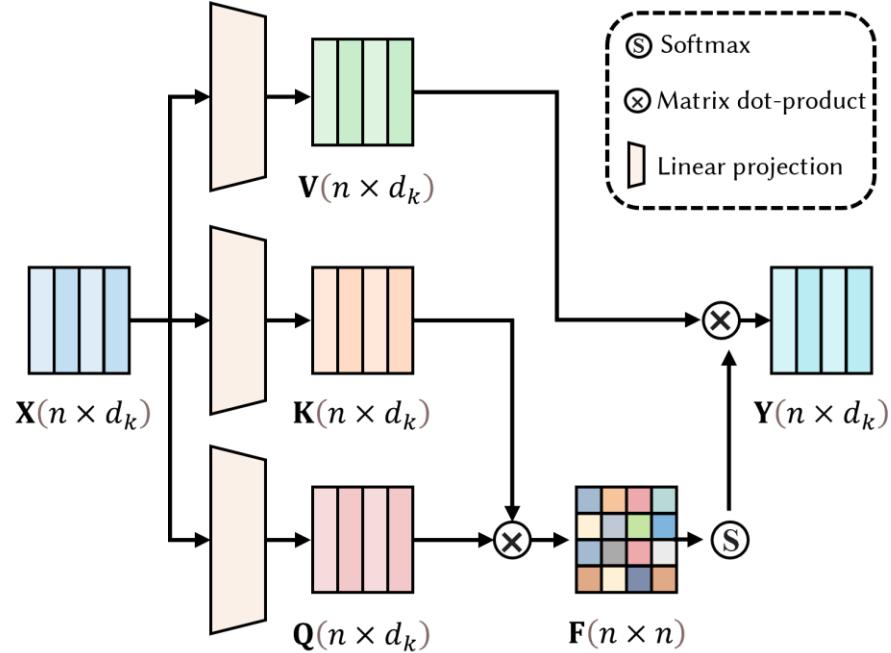
Challenge:

1. Variable length
2. Fusing timestep
3. Improve efficiency

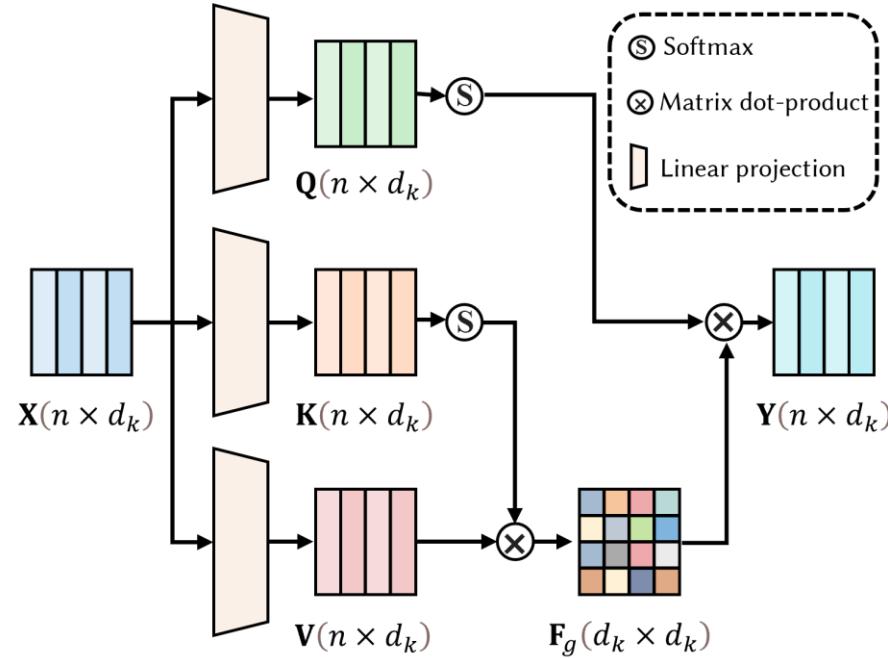
Cross-Modality Linear Transformer



Linear Self-Attention



Classical Self-Attention

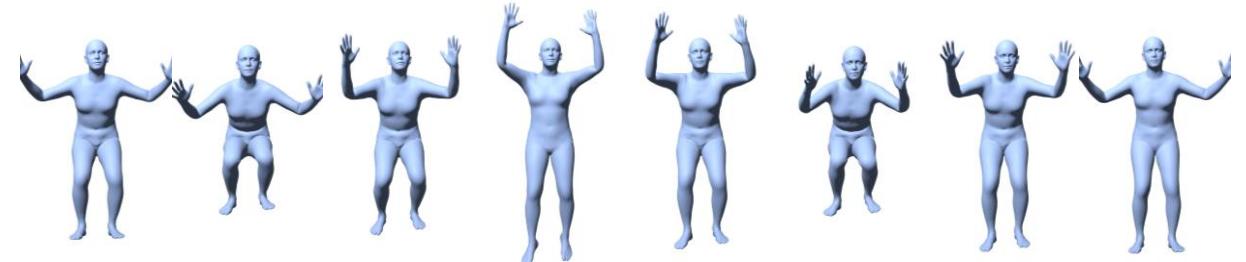


Linear Self-Attention

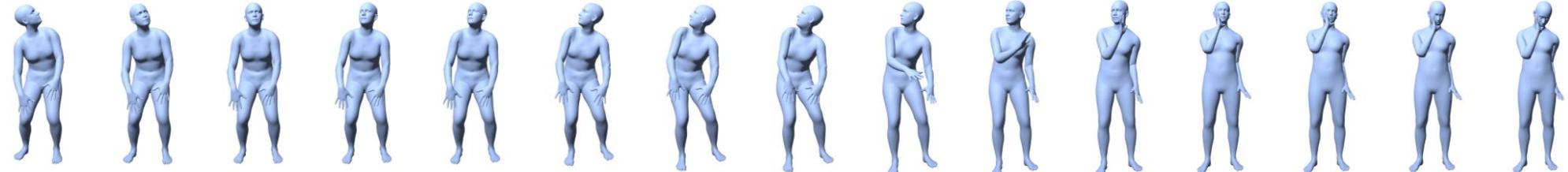
Visualization



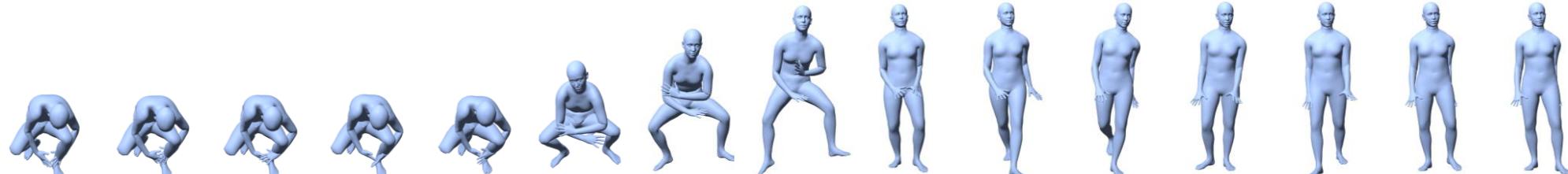
a) *Kicking and punching*



b) *Jumping and raising arms*



c) *Looking around and then calling on a phone with right hand*

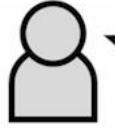


d) *Tying the shoe, standing up and then walking forward*





play the guitar



play the piano



play the violin



Examples

prompt

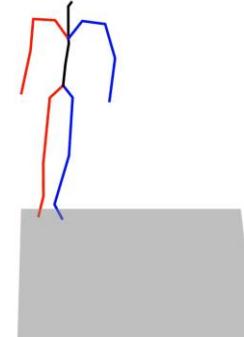
a person spins quickly and takes off running.

length

29

Clear **Submit**

a person spins quickly and takes off running. #29



prompt

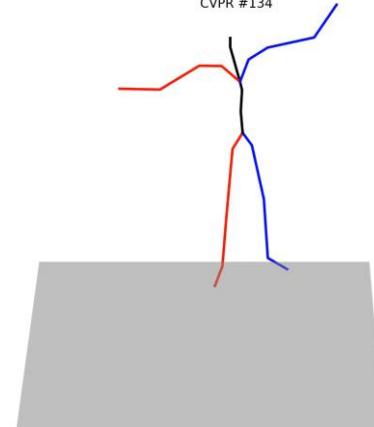
the man turning around and happy about the acceptance of CVPR

length

134

Clear **Submit**

the man turning around and happy about the acceptance of
CVPR #134



Online Demo

MotionDiffuse: Text-Driven Human Motion Generation with Diffusion Model

This is an interactive demo for MotionDiffuse. For more information, feel free to visit our project page(<https://mingyuan-zhang.github.io/projects/MotionDiffuse.html>).

prompt

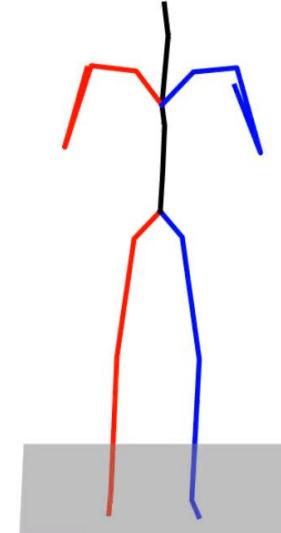
length

58

Clear Submit

output

the man throws a punch with each hand. #58



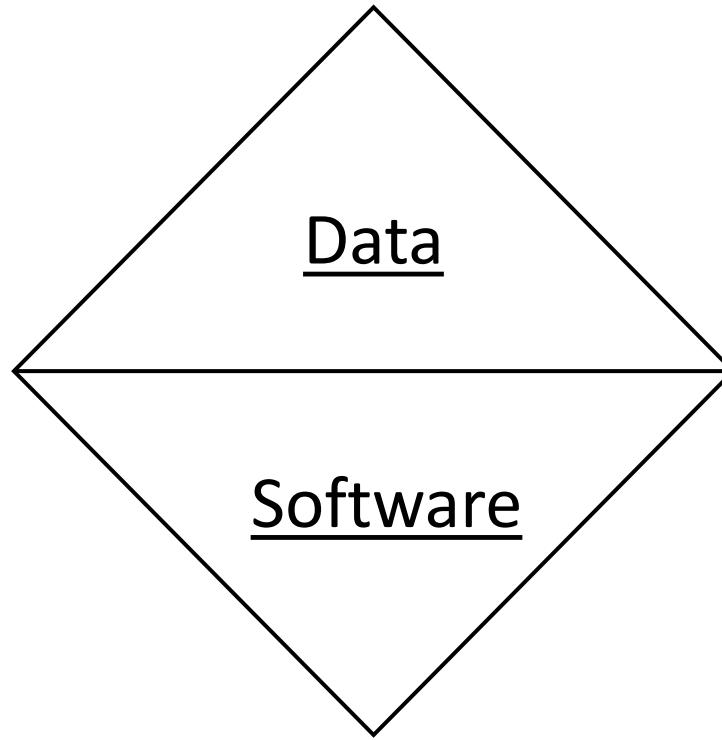




2D Generation



Motion Generation



3D Generation



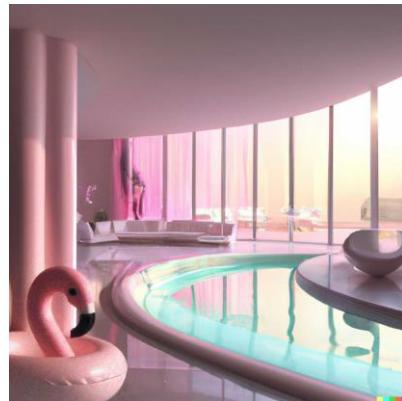
Scene Generation



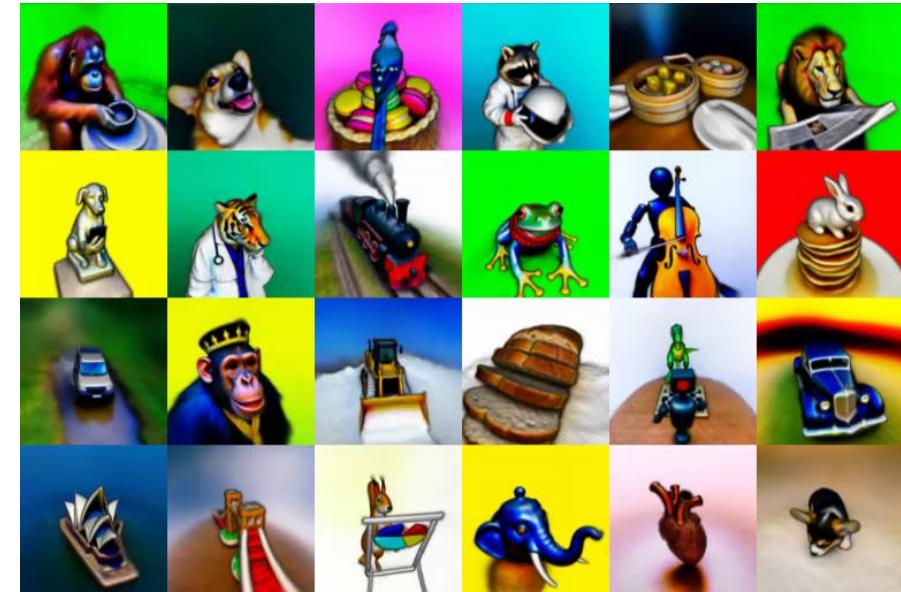
Text-driven Content Generation



Imagen^[1]



DALL·E^[2]



DreamFusion^[3]

[1] <https://imagen.research.google/>

[2] <https://openai.com/dall-e-2/>

[3] <https://dreamfusion3d.github.io/>

What about creating the environment?



The surrounding environment is also important to
an immersive VR experience.

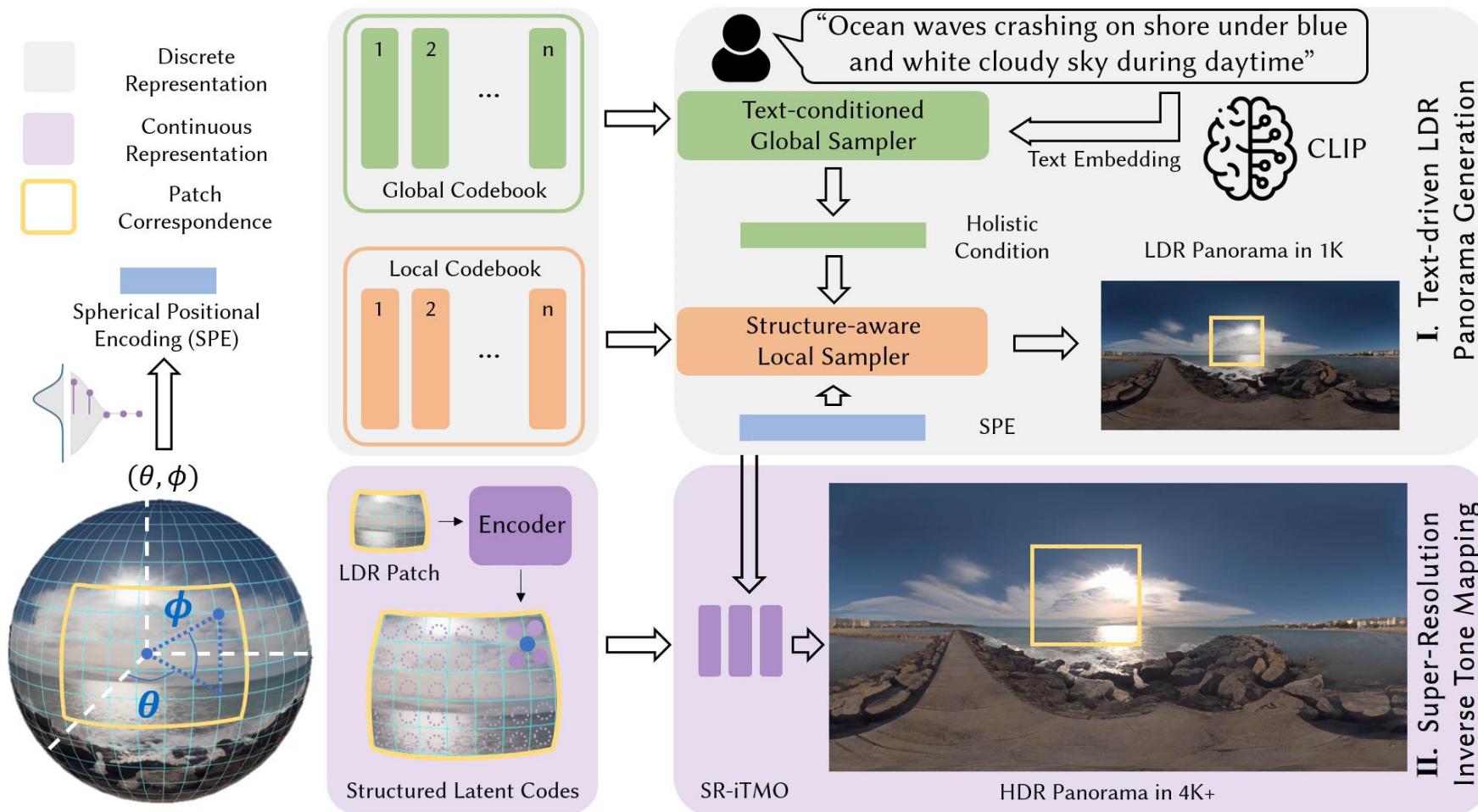


Full field of view (360°) → Panorama
Realistic illuminations → HDR
High-quality textures → 4K resolution

Create the Surroundings Using Texts

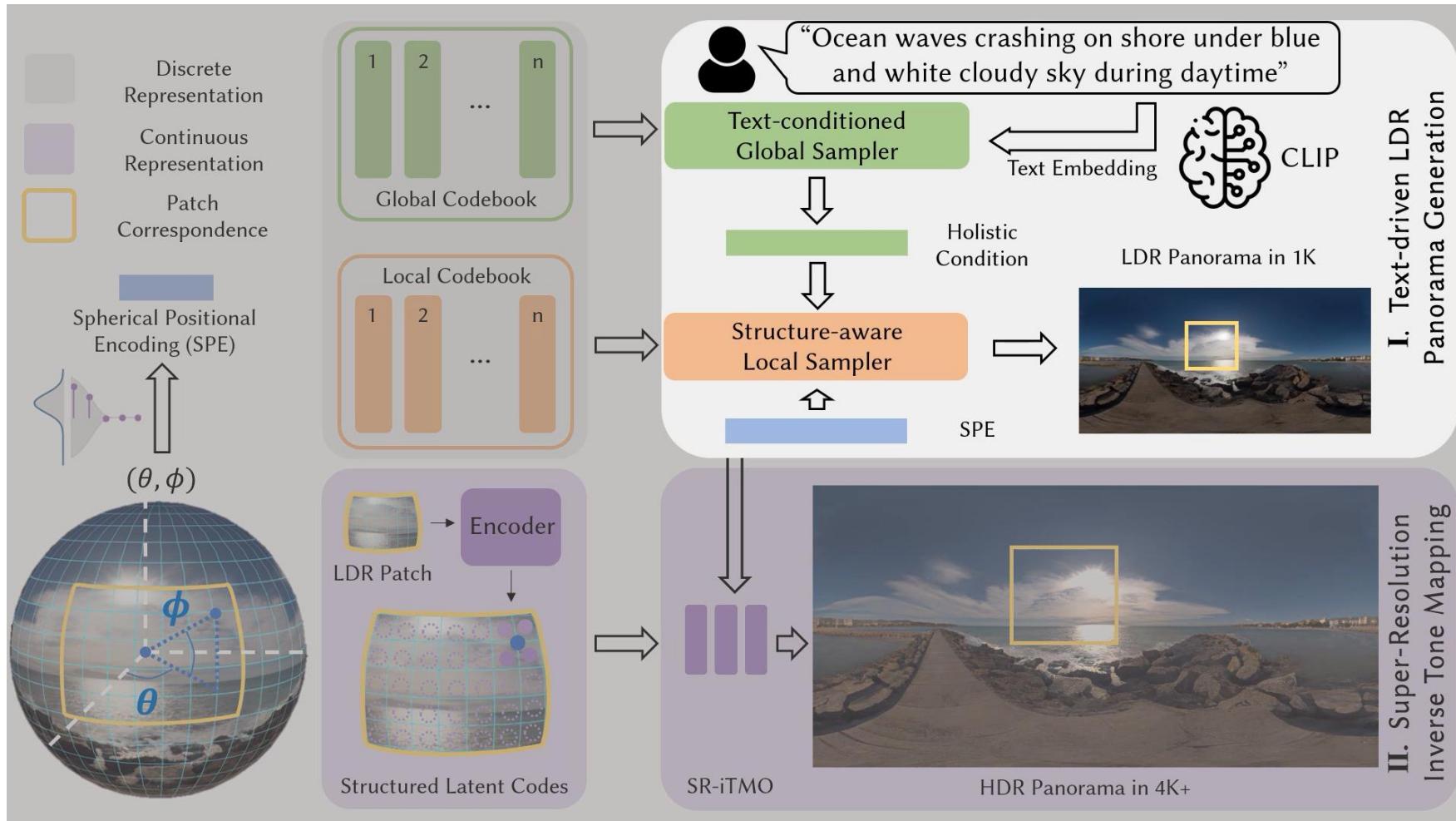


Text2Light An Overview



Text2Light

Stage I: Text-driven LDR Panorama Generation



Text2Light

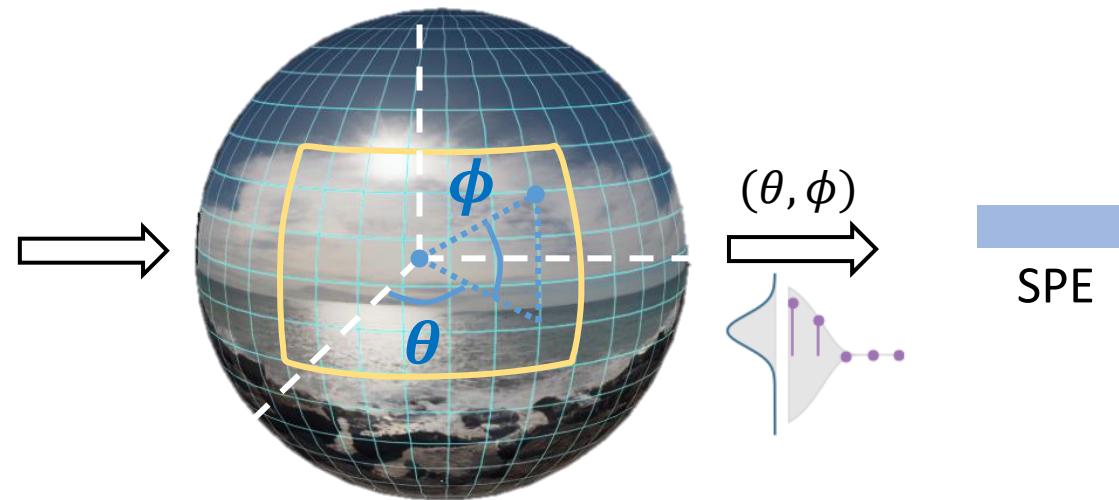
Stage I: Structure-aware Local Sampler



Spherical Positional Encoding (SPE)



Raw HDR Panorama

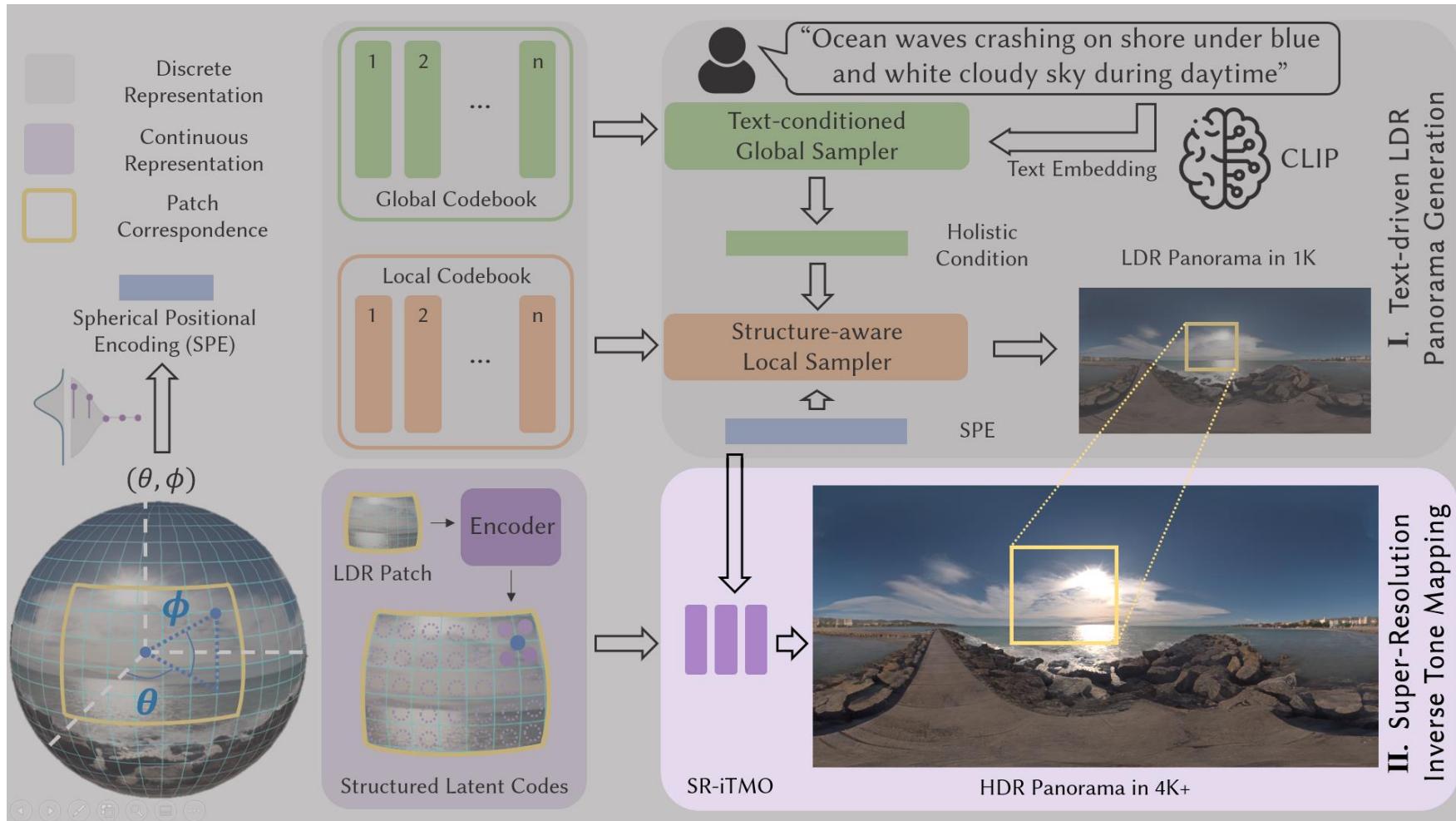


Spherical Fields

Fourier Features

Text2Light

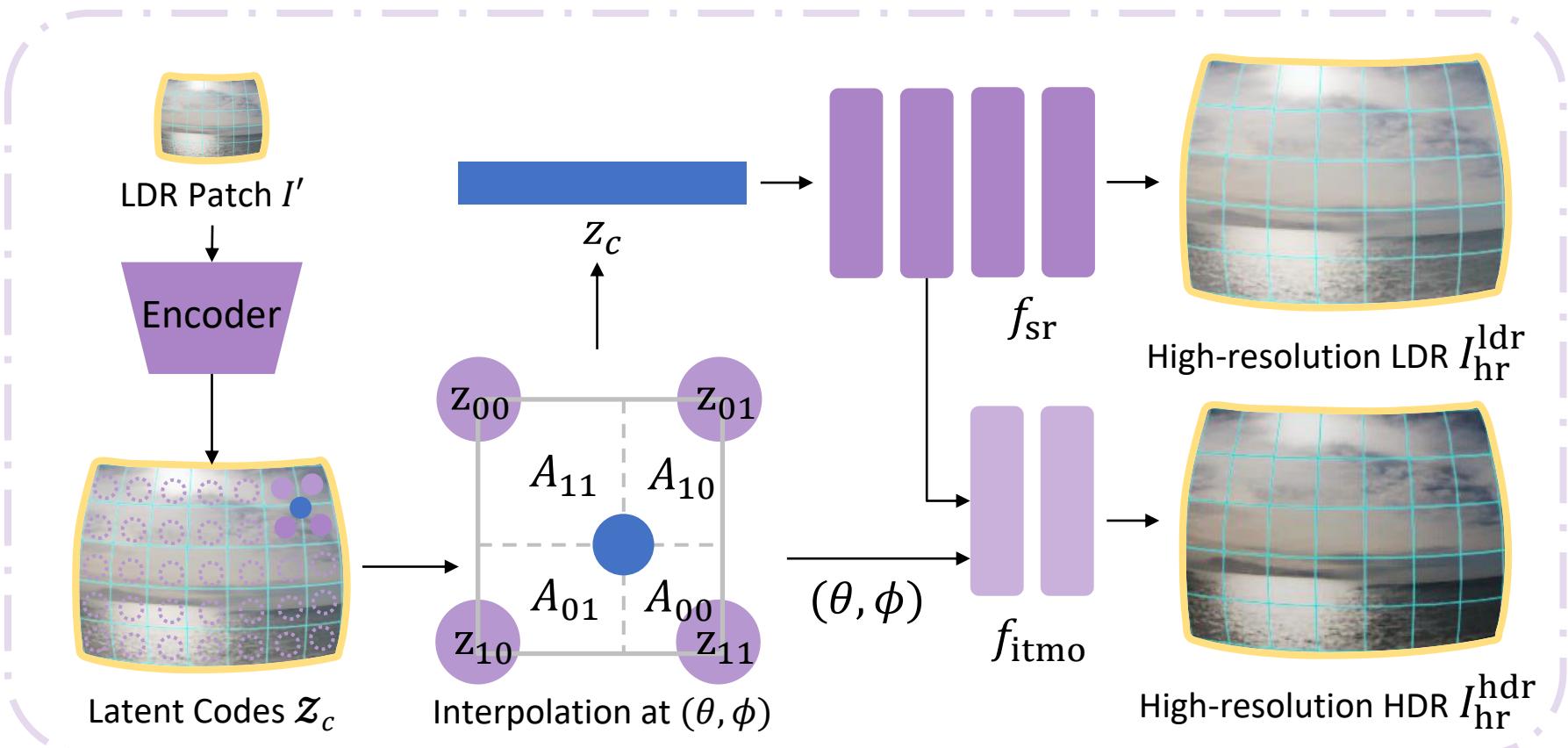
Stage II: Super-Resolution Inverse Tonemapping



Text2Light

Stage II: SR-iTMO as two MLPs

Super-Resolution Inversed Tone Mapping Operator (SR-iTMO)



LDR – Low Dynamic Range
HDR – High Dynamic Range

LR – Low Resolution
HR – High Resolution

Text2Light Applications: UI



The image shows a user interface for the Text2Light application. On the left, there is a large, semi-transparent green circular overlay. To the right, the application's branding is displayed: "Text2Light" in a stylized font where "Text" is green and "Light" is orange, with the tagline "Own Your Reality with Any Sentences" underneath. Below this, there is a text input field labeled "Describe Your Scene" with the placeholder "e.g. a living room". At the bottom, there are two green buttons labeled "Generate" and "Render".



“white bed
linen with
white pillow”



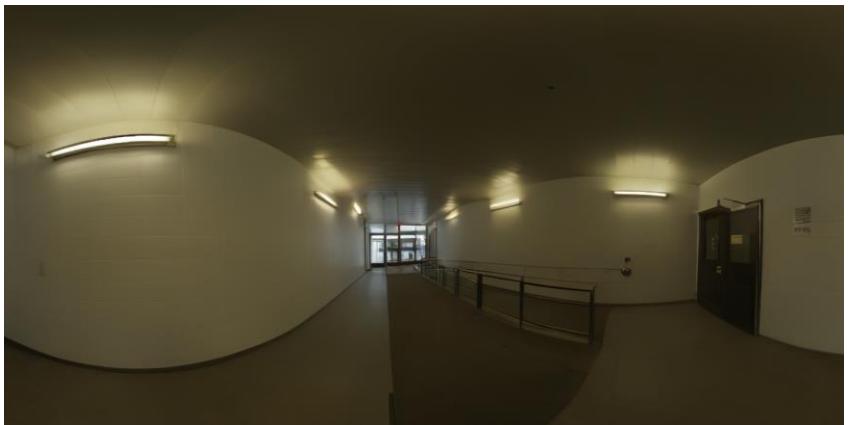
“brown wooden
floor with white
wall”



“gray concrete
pathway with
wall signages”



“closeup photo of
concrete stair
surrounded by
white painted wall”



“blue and
brown wooden
counter”



“empty parking
lot during
daytime”



Suzanne Monkey: glossy Shader balls: glass, diffuse, glossy, mixture of diffuse and glossy

“white bed
linen with
white pillow”



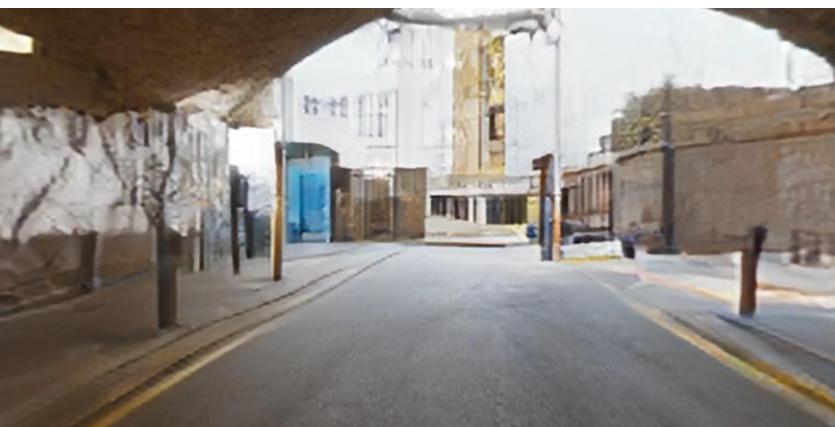
“gray concrete
pathway with
wall signages”



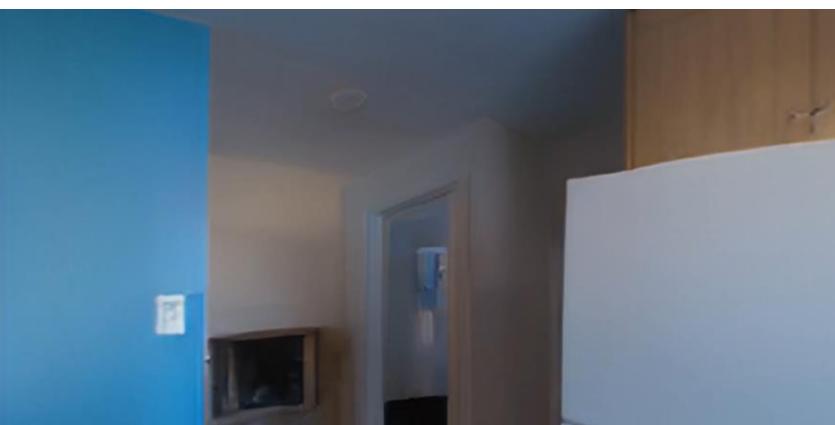
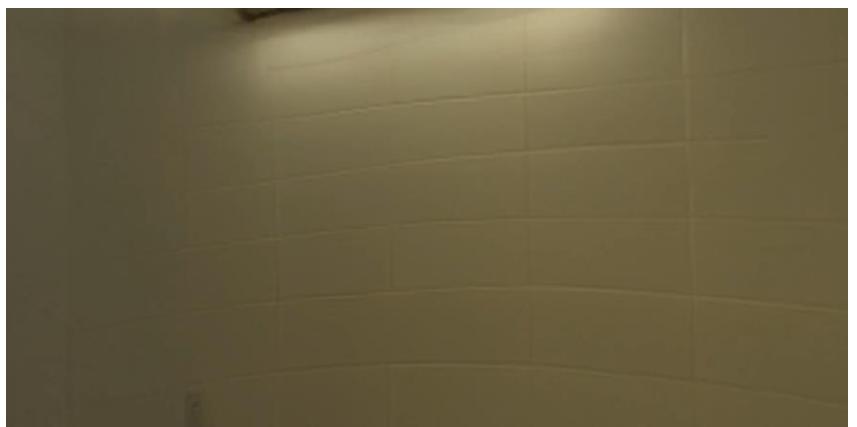
“blue and
brown wooden
counter”



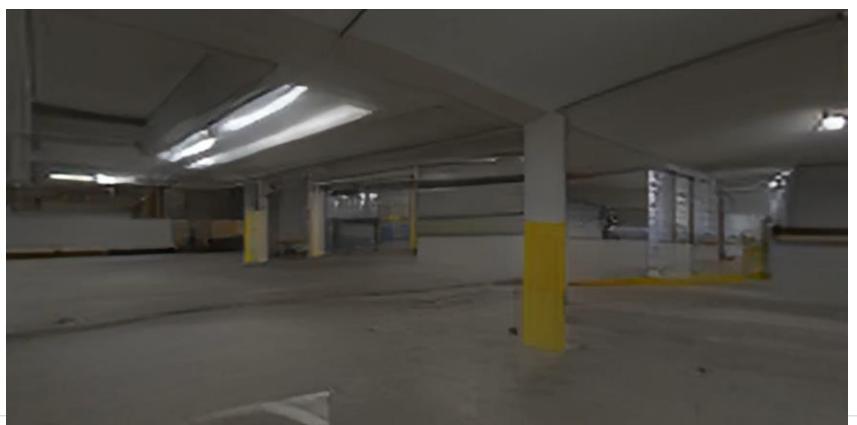
“brown wooden
floor with white
wall”



“closeup photo of
concrete stair
surrounded by
white painted wall”



“empty parking
lot during
daytime”



Suzanne Monkey: glossy Shader balls: glass, diffuse, glossy, mixture of diffuse and glossy

"lined brown
pew benches"



"Avenue, Trees,
Path, Sunbeams,
Sunrays"



"photo of
orange chairs"



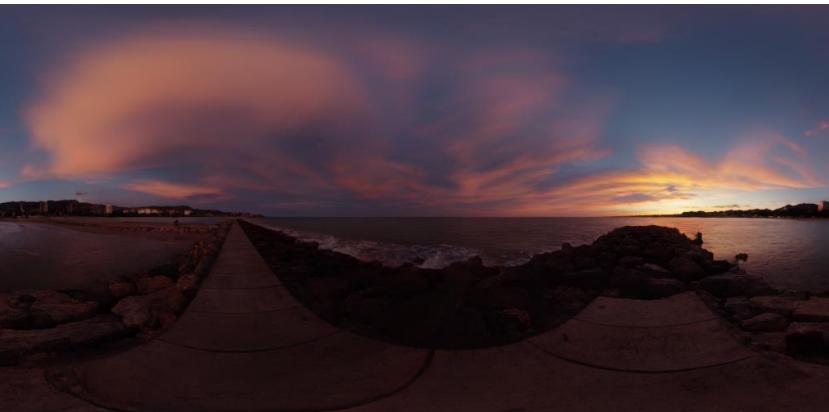
"ocean waves crashing
on shore under blue
and white cloudy sky
during daytime"



"road with falling
leaves in between
of trees"



"sunset by the
ocean"



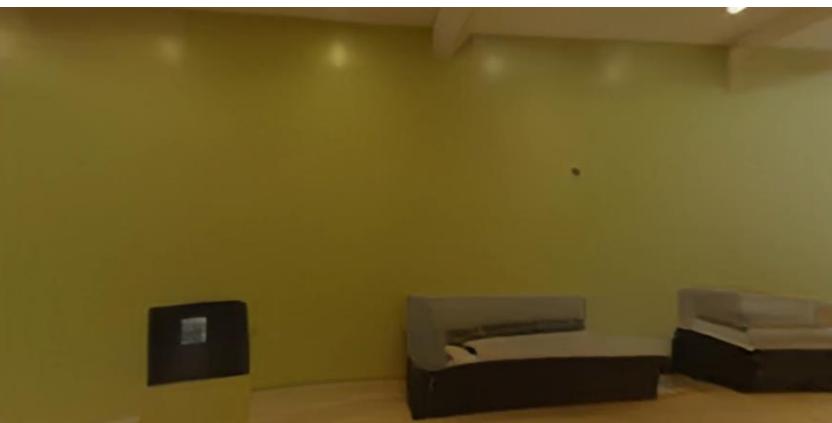
"lined brown
pew benches"



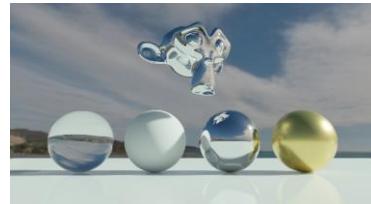
"Avenue, Trees,
Path, Sunbeams,
Sunrays"



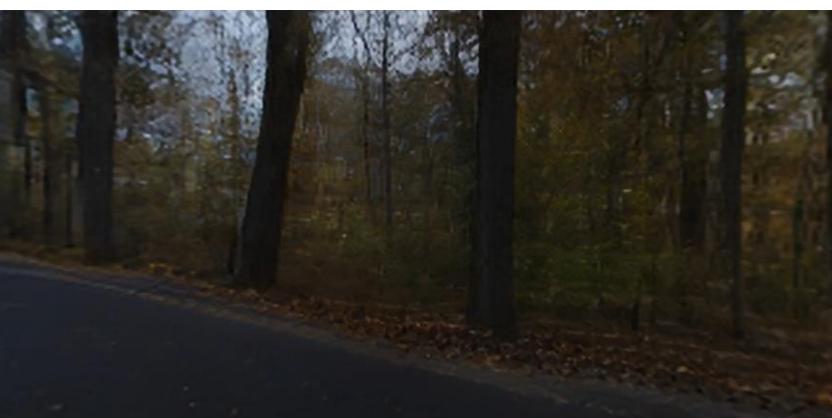
"photo of
orange chairs"



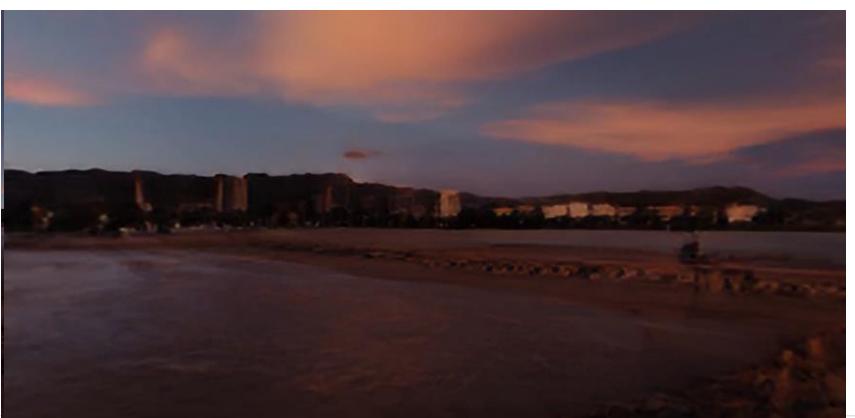
"ocean waves crashing
on shore under blue
and white cloudy sky
during daytime"



"road with falling
leaves in between
of trees"

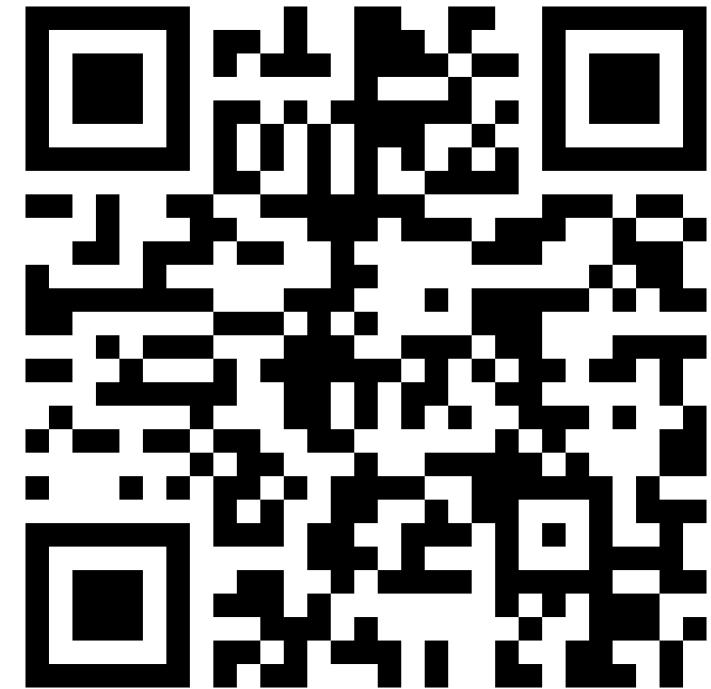


"sunset by the
ocean"



Suzanne Monkey: glossy Shader balls: glass, diffuse, glossy, mixture of diffuse and glossy

Text2Light



Project Page

<https://frozenburning.github.io/projects/text2light/>

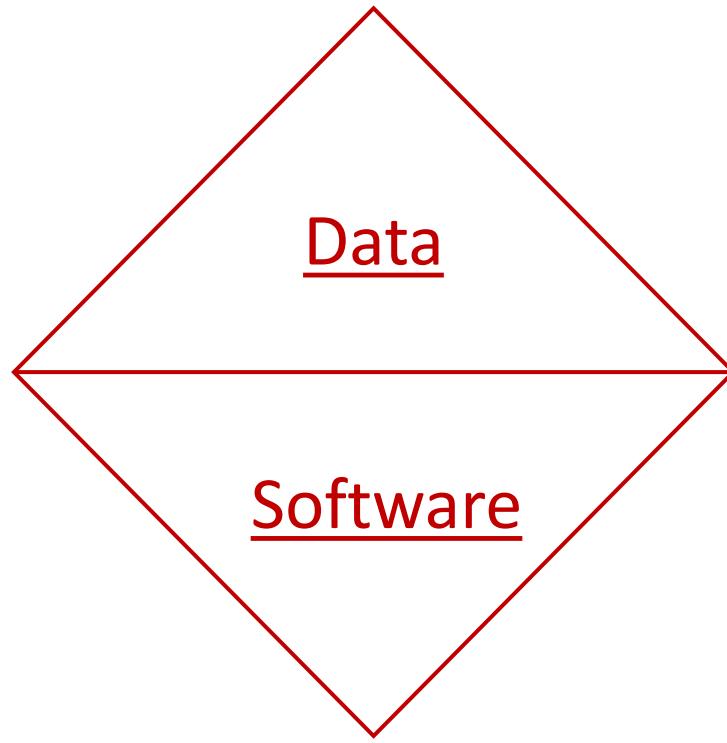




2D Generation



Motion Generation



3D Generation



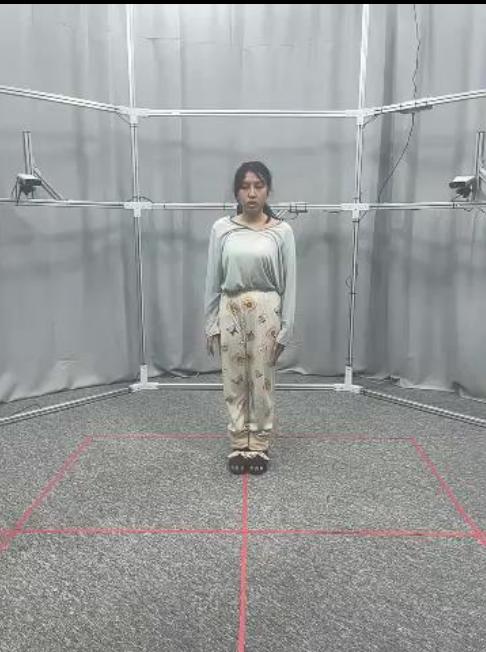
Scene Generation



HuMMan MoCap System



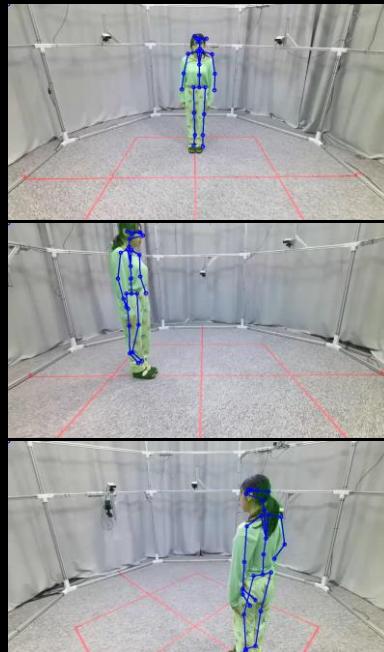
Artec Eva



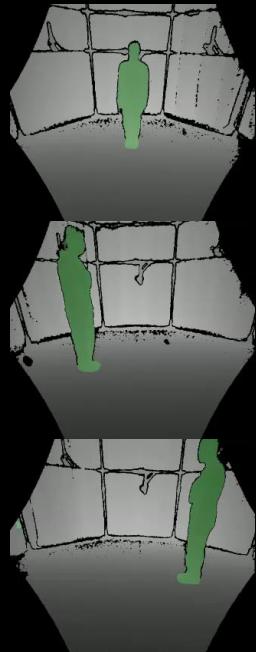
iPhone RGB



iPhone Depth



Kinect RGB



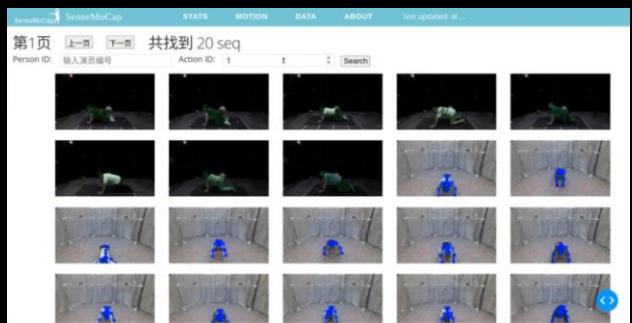
Kinect Depth

0.1mm
Accuracy

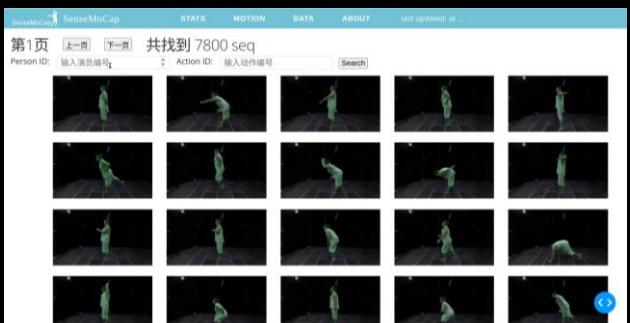
11
Cameras

1G
Data / Sec

6
Actor / Day



Search by Action



Search by Actor

Institute	#Camera	Mobile
Tencent	3	-
Tsinghua	6	-
ByteDance	6	-
NJU	9	-
CMU	10	-
Ours	≥ 10	✓

MMHuman3D Software





2D Generation



Motion Generation

Thank you!



3D Generation

 “brown wooden dock on lake surrounded by green trees during daytime”



Scene Generation

