# Multi-Modal Generative AI with Foundation Models

**Ziwei Liu    刘子纬**

**Nanyang Technological University**

**2023**

By ~~2027~~, creators won't have to be technical, just creative, thanks to automation tools.
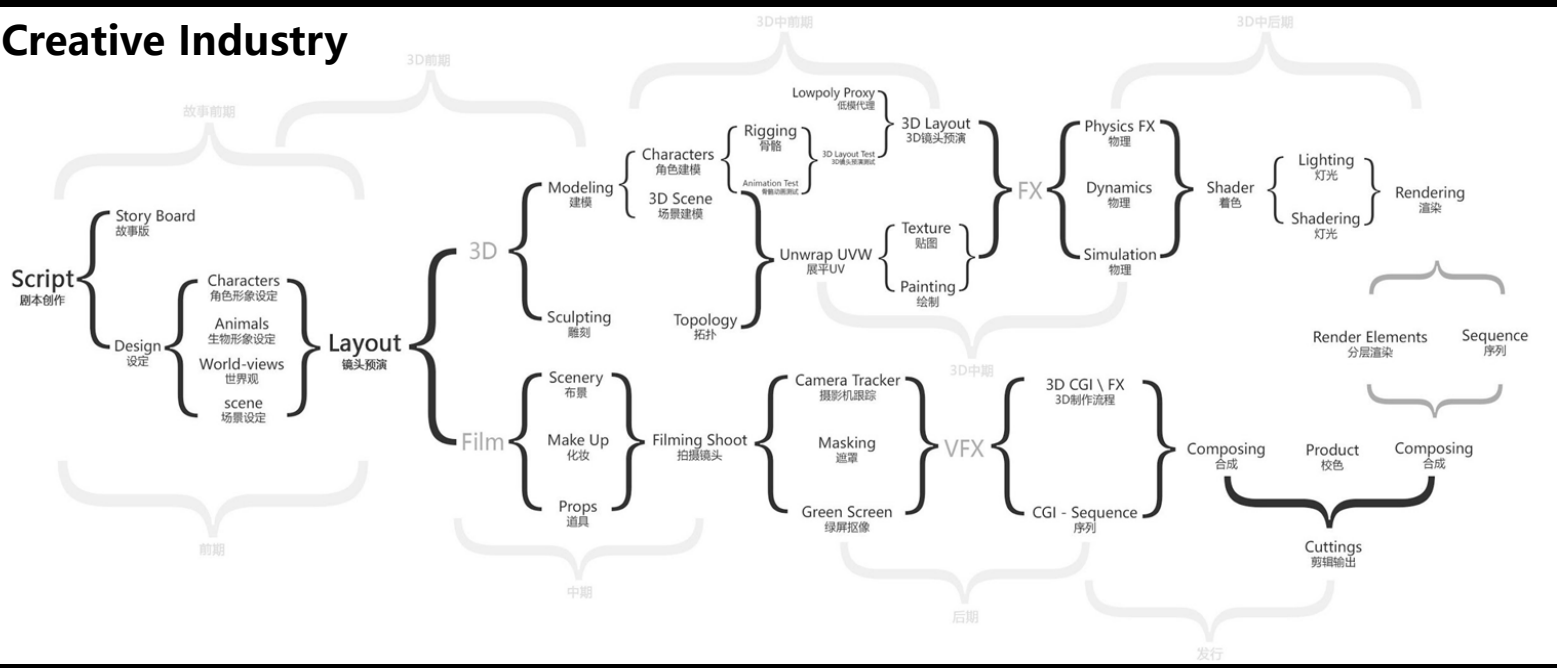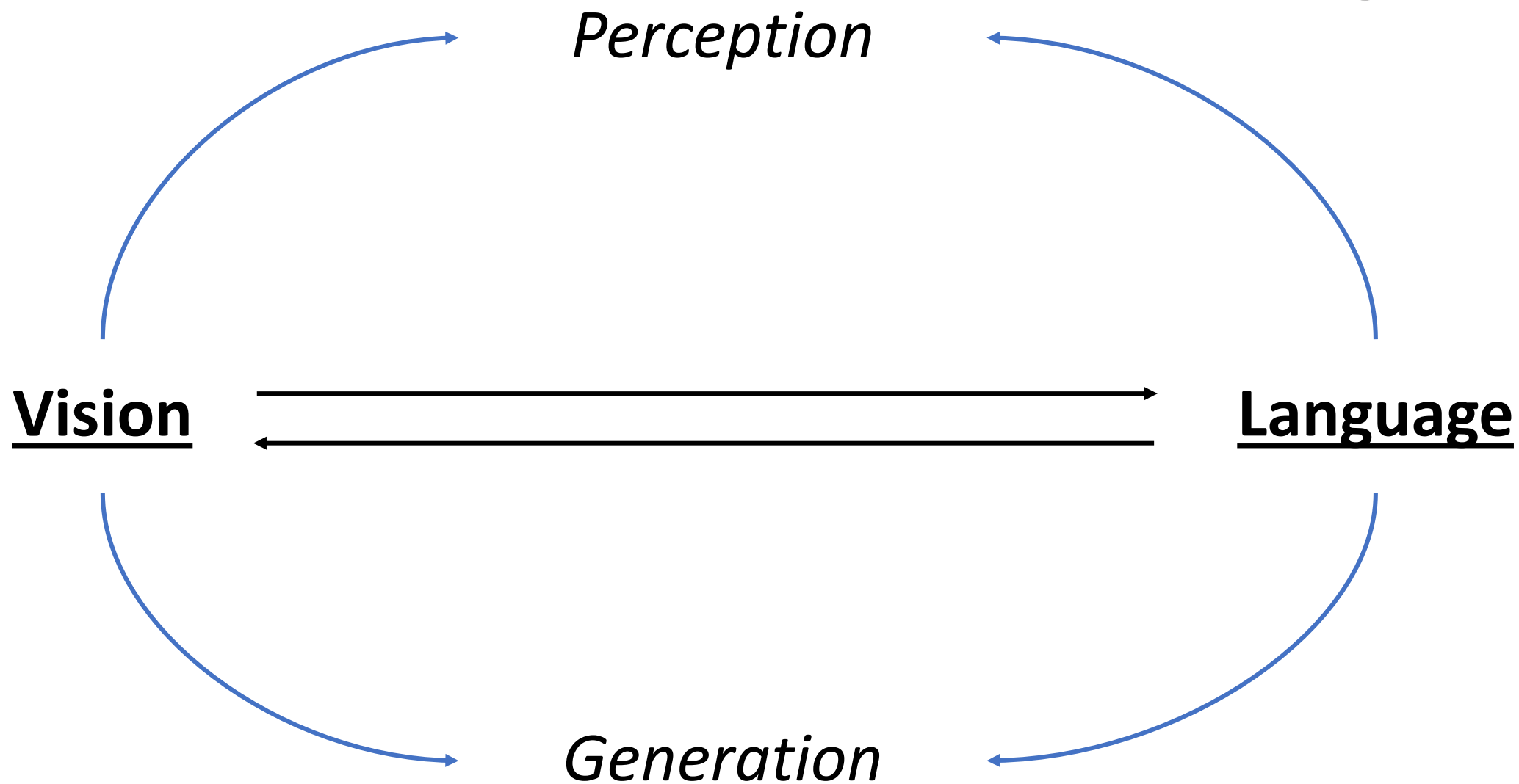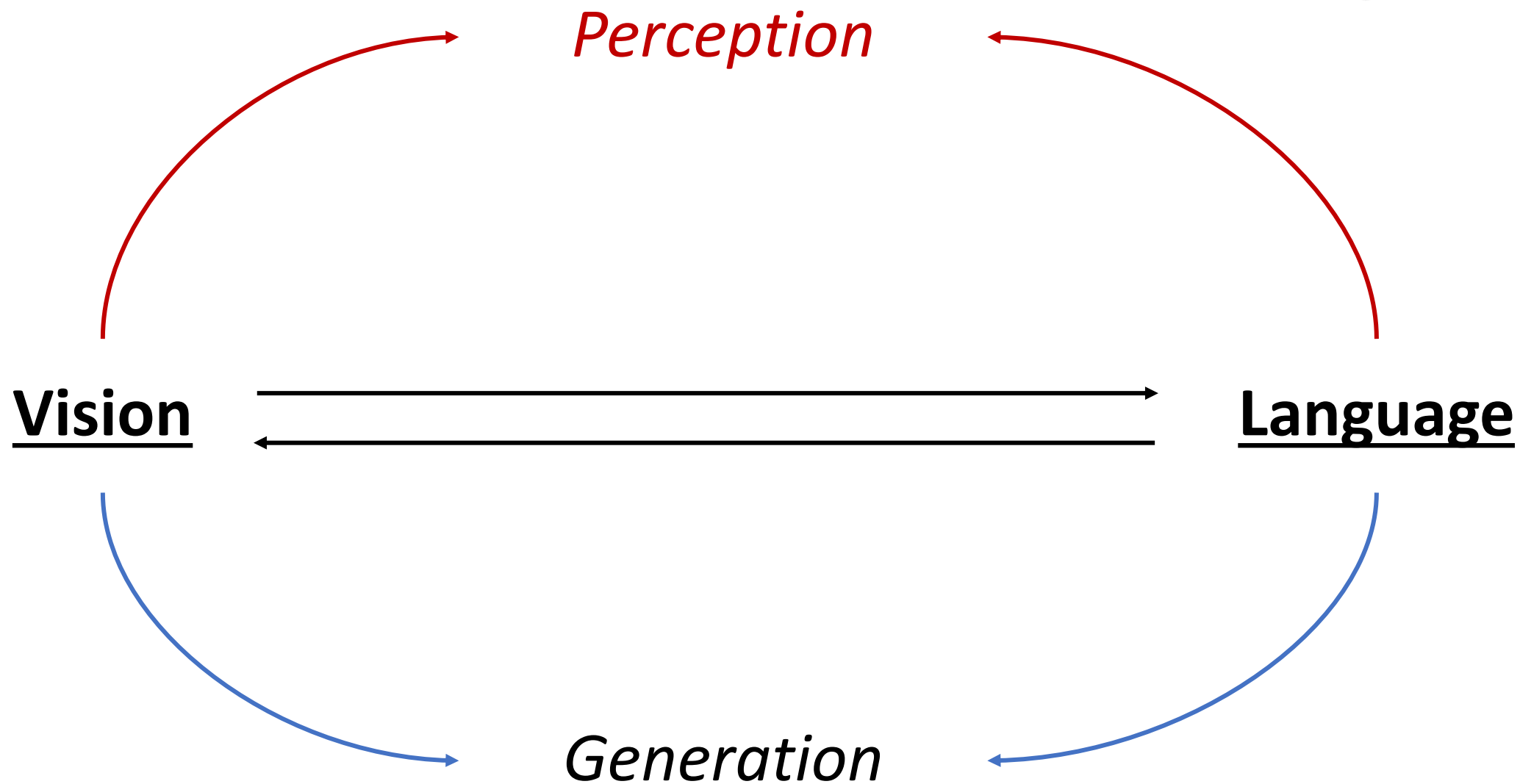
# AI-Generated Content


**Movie**


**Game**


**Anime**


**VTuber**


**Virtual Beings**

**Creative Industry**

# Training, Deployment and Evaluation of Foundation Models

# The pathway: From Language Models to Language Assistant

GPT-2    GPT-3    GPT-3.5    ChatGPT

Industrial

Open-source

BERT    LLaMA/T5    Vicuna/Flan-T5    Open Assistant

Zero-shot learning

Zero-shot learning
In-context learning

Zero-shot learning
In-context learning
Instruct following

Zero-shot learning
In-context learning
Instruct following
Human alignment

# The pathway: From Multi-modal Models to Multi-modal Assistants

CLIP

Flamingo

OpenAI

Google DeepMind

Industrial

Open-source

OpenCLIP

OpenFlamingo

LAION

LAION

Zero-shot learning

Zero-shot learning
In-context learning

# The pathway: From Multi-modal Models to Multi-modal Assistants

CLIP

Flamingo

OpenAI

Google DeepMind

Industrial

Open-source

OpenCLIP

OpenFlamingo

LAION

LAION

Otter

Zero-shot learning

Zero-shot learning
In-context learning

# Flamingo: a Visual Language Model for Few-Shot Learning



Alayrac et. al. Flamingo: a visual language model for few-shot learning. 2022

# Perceiver: versatile to multiple images and in-context examples



Input webpage → Processed text: <image> tags are inserted and special tokens are added
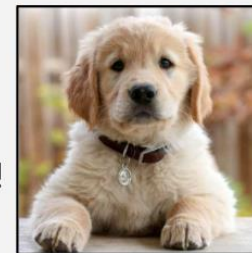
Image-Text Pairs dataset
[N=1, T=1, H, W, C]

Video-Text Pairs dataset
[N=1, T>1, H, W, C]

Multi-Modal Massive Web (M3W) dataset
[N>1, T=1, H, W, C]

# Flamingo Application

## Zero-shot learning



Input Prompt

Question: What do you think the capacities of these are? Answer:

Completion

**The floppy disk is 1.44MB and the CD is 700MB.**

Input Prompt

Question: What nutrients is the fruit in the image rich in? Answer:

Completion

**Vitamin C, potassium, and fiber.**

## In-context learning

### Input Prompt

This is a chinchilla. They are mainly found in Chile.

This is a shiba. They are very popular in Japan.

This is

### Completion

**a flamingo. They are found in the Caribbean and South America.**

## Video Understanding

What happens to the man after hitting the ball? Answer:

**he falls down.**

# Flamingo Application

**multi-image visual dialogue**

# Flamingo ≠ Multi-modal Assistants

OpenFlamingo simply completes the next reasonable sentence.

...he danger of this sport?

**OpenFlamingo\***: What is the danger of playing baseball? What is the danger of this sport? What might be the danger of this sports?

Flaming (trained in the SSL manner) are not aligned with user intent and serve as a Chatbot.

*OpenFlamingo is the open-source version of Flamingo, enabling community research with a strong interleaved data pretrained model

# Flamingo ≠ Multi-modal Assistants



Question: What is the danger of this sport?

**Human Expected**: The sport involves players running and trying to catch the ball while others are standing in the grass, which can lead to collisions or accidents.

Flaming (trained in the SSL manner) are not aligned with user intent and serve as a Chatbot.
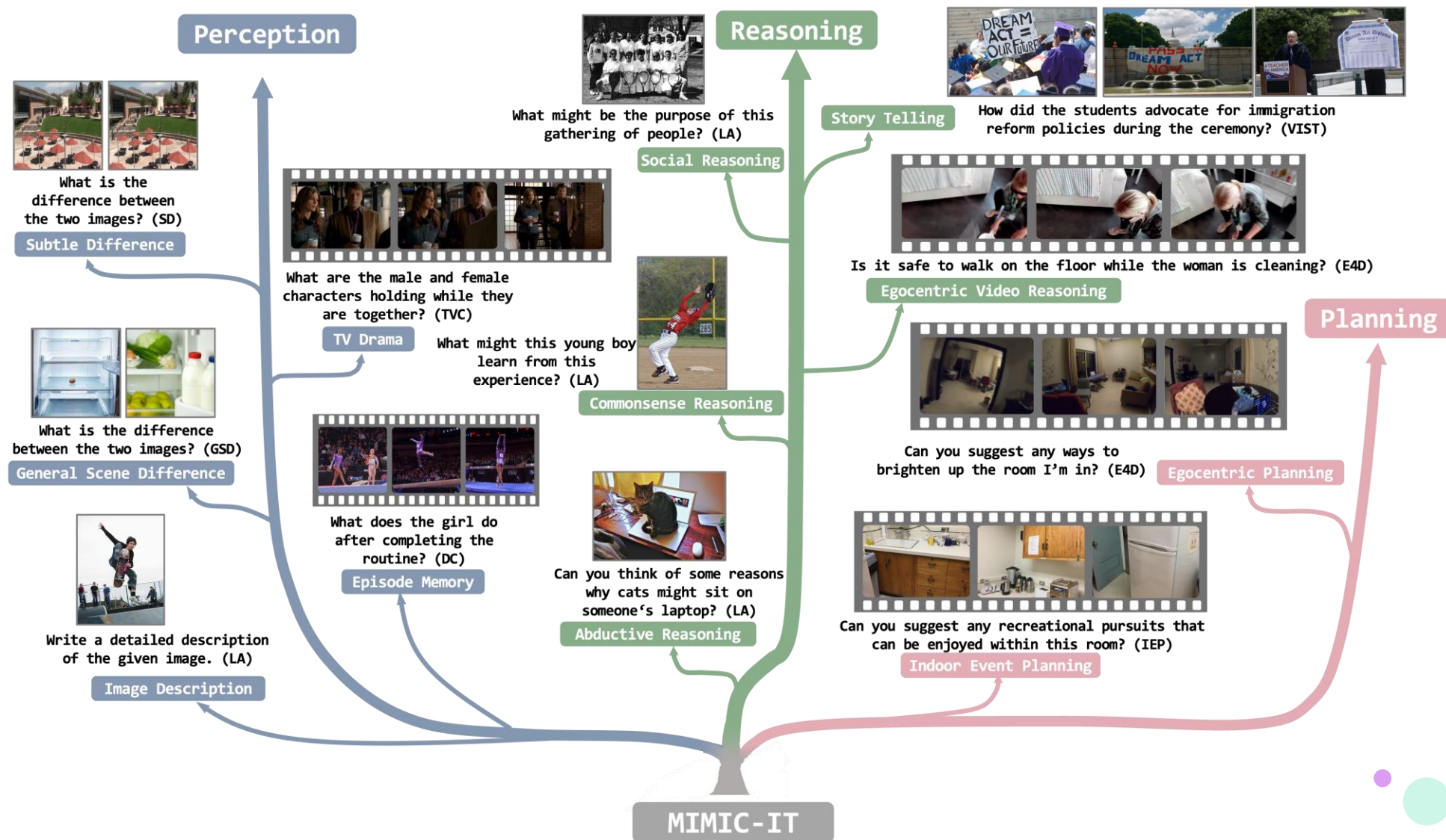
# MMC4: Image-text interleaved data for OpenFlamingo Pretraining



[..., "Check out Shane Driscoll's take on sustainable communities and how his photograph fits this year's Green Cities theme.", ..., ,"Man-made platforms like the one pictured here allow these fish-eating birds of prey to thrive in developed coastal areas.", "A city surrounded by mountains.", "I took this photo in October on a hike in New Hampshire.", , "It is looking at Mt. Chicora from the middle sister mountain.", "Getting people out into beautiful places like this is becoming more and more popular, and each time we bring a little piece of nature back with us that inspires us to make our cities better.", ...]
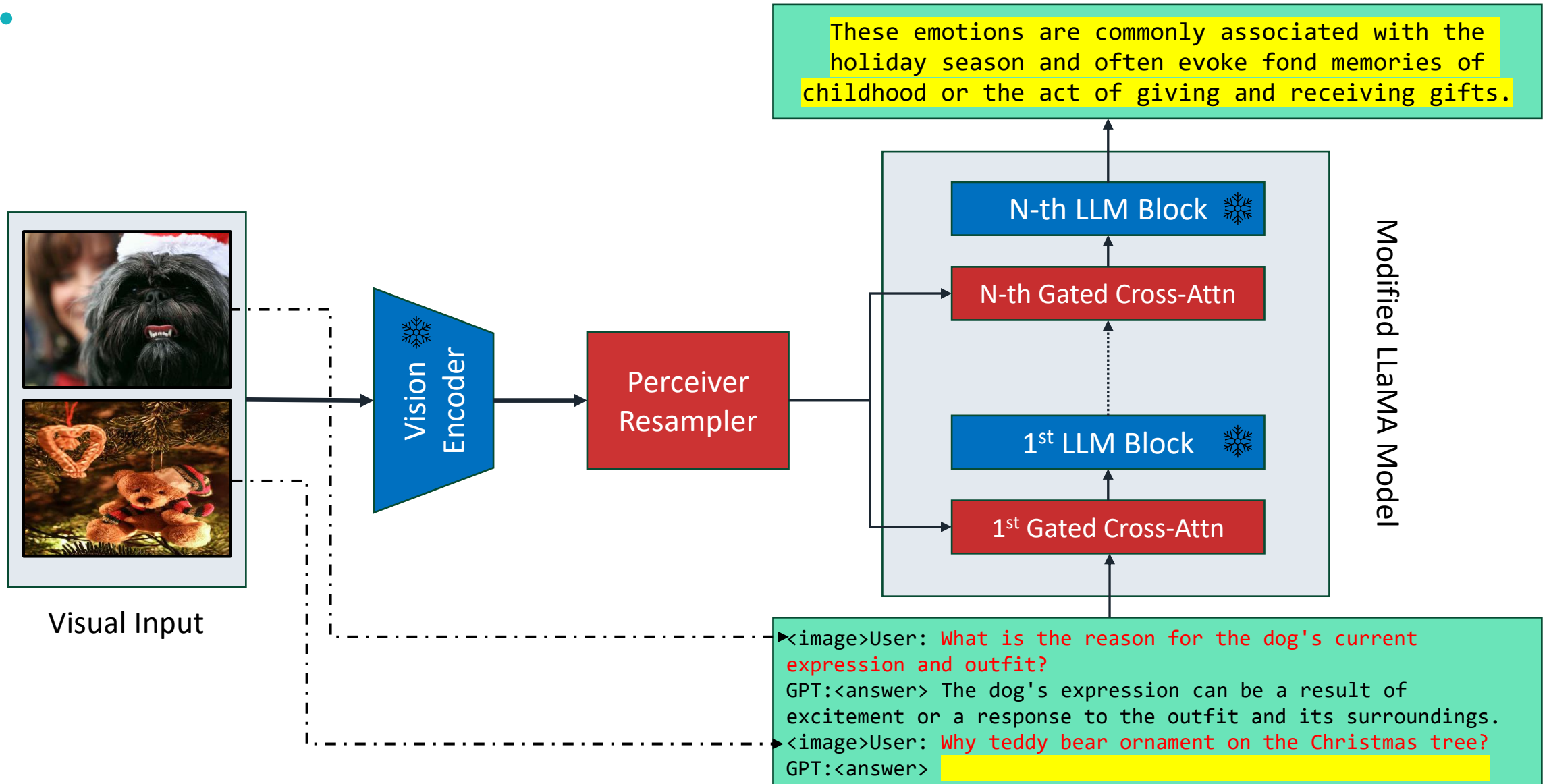
Diverse and large-scale, but lack of Instruct-following scenario

Zhu et. al. Multimodal C4: An open, billion-scale corpus of images interleaved with text. 2023

# MIMIC-IT Dataset



**Perception**

What is the difference between the two images? (SD)
*Subtle Difference*

What is the difference between the two images? (GSD)
*General Scene Difference*

Write a detailed description of the given image. (LA)
*Image Description*

What are the male and female characters holding while they are together? (TVC)
*TV Drama*

What does the girl do after completing the routine? (DC)
*Episode Memory*

**Reasoning**

What might be the purpose of this gathering of people? (LA)
*Social Reasoning*

What might this young boy learn from this experience? (LA)
*Commonsense Reasoning*

Can you think of some reasons why cats might sit on someone's laptop? (LA)
*Abductive Reasoning*

*Story Telling*

How did the students advocate for immigration reform policies during the ceremony? (VIST)

Is it safe to walk on the floor while the woman is cleaning? (E4D)
*Egocentric Video Reasoning*

**Planning**

Can you suggest any ways to brighten up the room I'm in? (E4D)
*Egocentric Planning*

Can you suggest any recreational pursuits that can be enjoyed within this room? (IEP)
*Indoor Event Planning*

**MIMIC-IT**

# Otter: A Multi-Modal In-context Instruction Tuned Model

# From interleaved data pretraining to multi-modal In-context instruction tuning



**MMC4**
(interleaved pretraining)

**OpenFlamingo**

**MIMIC-IT**
(Multi-Modal In-Context Instruction Tuning)

**Otter**

# Otter

## Cognition

Sum of the scores of all cognition subtasks, including commonsense reasoning, numerical calculation, text translation, and code reasoning. The full score of each subtask is 200, and that of all cognition is 800.

| Rank | Model | Version | Score |
|------|-------|---------|-------|
| 🥇 | **Otter** | **OTTER-Image-MPT7B** | **306.43** |
| 🥈 | **MiniGPT-4** | **minigpt4-aligned-with-vicuna13b** | **292.14** |
| 🥉 | **InstructBLIP** | **blip2-instruct-flant5xxl** | **291.79** |
| 4 | BLIP-2 | blip2-pretrain-flant5xxl | 290.00 |
| 5 | mPLUG-Owl | mplug-owl-llama-7b | 276.07 |
| 6 | LaVIN | LAVIN-13B | 249.64 |
| 7 | LLaMA-Adapter V2 | LLaMAv2-7B | 248.93 |
| 8 | PandaGPT | pandagpt-7b-max-len-512 | 228.57 |
| 9 | Multimodal-GPT | Multimodal-GPT-9B | 226.79 |
| 10 | LLaVA | LLaVA-7B-v0 | 214.64 |
| 11 | ImageBind_LLM | imagebind_LLM-7B | 213.57 |
| 12 | VisualGLM-6B | VisualGLM-6B | 181.79 |

Otter

Multi-Modal In-Context Learning
Model with Instruction Tuning

NANYANG TECHNOLOGICAL UNIVERSITY SINGAPORE

S-LAB FOR ADVANCED INTELLIGENCE

Multi-Modal In-Context Learning
Model with Instruction Tuning

NANYANG
TECHNOLOGICAL
UNIVERSITY
SINGAPORE

S-LAB
FOR ADVANCED
INTELLIGENCE

# Contrastive Language-Image Pre-training (CLIP)

- Data: 400M image-text pairs
- Compute: 250-600 GPUs
- Train time: up to 18 days



Radford et al. Learning Transferable Visual Models From Natural Language Supervision. ICML'21.

# Zero-shot image recognition via prompting



Radford et al. Learning Transferable Visual Models From Natural Language Supervision. ICML'21.

# Fine-tuning might not be a good idea



- Fine-tuning the image encoder: accuracy drops by ~40%
- Fine-tuning both encoders could lead to collapse

Zhou, Yang, Loy, and Liu. Learning to Prompt for Vision-Language Models. IJCV'22.

# Prompt engineering is too time-consuming

| Caltech101 | Prompt | Accuracy |
|---|---|---|
|  | a [CLASS]. | 82.68 |
| | a photo of [CLASS]. | 80.81 |
| | a photo of a [CLASS]. | 86.29 |
| | $[V]_1 [V]_2 \dots [V]_M$ [CLASS]. | **91.83** |

| Flowers102 | Prompt | Accuracy |
|---|---|---|
|  | a photo of a [CLASS]. | 60.86 |
| | a flower photo of a [CLASS]. | 65.81 |
| | a photo of a [CLASS], a type of flower. | 66.14 |
| | $[V]_1 [V]_2 \dots [V]_M$ [CLASS]. | **94.51** |

| Describable Textures (DTD) | Prompt | Accuracy |
|---|---|---|
|  | a photo of a [CLASS]. | 39.83 |
| | a photo of a [CLASS] texture. | 40.25 |
| | [CLASS] texture. | 42.32 |
| | $[V]_1 [V]_2 \dots [V]_M$ [CLASS]. | **63.58** |

| EuroSAT | Prompt | Accuracy |
|---|---|---|
|  | a photo of a [CLASS]. | 24.17 |
| | a satellite photo of [CLASS]. | 37.46 |
| | a centered satellite photo of [CLASS]. | 37.56 |
| | $[V]_1 [V]_2 \dots [V]_M$ [CLASS]. | **83.53** |

Zhou, Yang, Loy, and Liu. Learning to Prompt for Vision-Language Models. IJCV'22.

# Context Optimization (CoOp)



*Key insight: update input prompt parameters*

Zhou, Yang, Loy, and Liu. Learning to Prompt for Vision-Language Models. IJCV'22.

# CoOp is a few-shot learner



Zhou, Yang, Loy, and Liu. Learning to Prompt for Vision-Language Models. IJCV'22.

# NOAH



(a)

(b)

# Visual In-Context Learning



(a) Visual in-context learning is sensitive to prompt selection

# Content Generation Powered by Foundation Models

# Visual AIGC

Avatar

Foundation
Model

Object

Scene

# Visual AIGC

**Avatar**

**Foundation Model**

Object

Scene

# StyleGAN-Human: 2D Human Generation

# Text2Human: Text-to-2D Human

# Text2Performer: Text-to-2D Human Video

The dress the person wears has medium sleeves and
it is of short length. The texture of it is pure color.

The lady moves to the left.

She is turning right from the front to the side.

She is turning right from the side to the back.

She turns right from the back to the side.

She turns right from the side to the front.

She moves to the right.

# EVA3D: 3D Human Generation

- Learn 3D generation from 2D image collections



*Static* ➡ **?** *Articulated*

# EVA3D: 3D Human Generation

- Compositional Human NeRF



**a)** Sample Rays & Points at Observation Space SMPL($\beta$, $\theta$)

**b)** Inverse LBS to Canonical Space

**c)** Transform to Local Bounding Box & Query Subnetworks

**d)** Integrate Queried Results Along Rays

# EVA3D: 3D Human Generation

- Qualitative Results

# EVA3D: 3D Human Generation

- Explicit Pose/ Shape Control

# AvatarCLIP: Text-to-3D Avatar

Cosine Distance

$E_I$

$E_T$

$t$ = "*Iron Man*"

CLIP Image Encoder

CLIP Text Encoder

a) Differentiable Rendering        b) Optimization guided by CLIP

# AVATARCLIP: DETAILED PIPELINE

$t_{shape}$ = "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

$t_{shape}$ = "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

CONTROLLING & CONCEPT MIXING ABILITIES

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

# AvatarCLIP: Text-to-3D Avatar

# MotionDiffuse: Text-to-3D Human Video

# 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**

# Motion Generation with Diffusion Model

**Diffusion Process**



$$\mathbf{x_0} \sim q(\mathbf{x_0})$$

Add Noise   Add Noise   Add Noise

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

**Reverse Process**



$$\mathbf{x_0} \sim q(\mathbf{x_0})$$

Denoise   Denoise   Denoise

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

# Framework



## Challenge：
1. Variable length
2. Fusing timestep
3. Improve efficiency

# Text-driven Motion Generation

# ReMoDiffuse: Text-to-3D Human Video

# HuMMan Dataset



**Artec Eva**  **iPhone RGB**  **iPhone Depth**  **Kinect RGB**  **Kinect Depth**

**Search by Action**  **Search by Actor**

**0.1mm** Accuracy

**11** Cameras

**1G** Data / Sec

**6** Actor / Day

# MMHuman3D Software


Input Video


Online Motion Capture



3D Animation Production:

3 days -> 30 min

# Visual AIGC



Avatar

Foundation Model

Object

Scene

# OmniObject3D: Text-to-3D Object

**OmniObject3D** is a **large-vocabulary** 3D dataset for **real-world scanned objects**.

- ✓ **6k** high-quality 3D models
- ✓ **190** categories
- ✓ **4** modalities: textured mesh, point cloud, real-captured video, synthetic multi-view images.
- ✓ Many down-stream tasks

| Dataset | Year | Real | Full 3D | Video | Num Objs | Num Cats |
|---------|------|------|---------|-------|----------|----------|
| ShapeNet | 2015 | | √ | | 51k | 55 |
| ModelNet | 2014 | | √ | | 12k | 40 |
| 3D-Future | 2020 | | √ | | 16k | 34 |
| ABO | 2021 | | √ | | 8k | 63 |
| Toys4K | 2021 | | √ | | 4k | 105 |
| CO3D | 2021 | √ | | √ | 19k | 50 |
| DTU | 2014 | √ | √ | | 124 | NA |
| GSO | 2021 | √ | √ | | 1k | 17 |
| AKB-48 | 2022 | √ | √ | | 2k | 48 |
| **Ours** | 2022 | √ | √ | √ | **6k** | **190** |

Real-world 3D scans

# Background and motivation

## Synthetic data

**ShapeNet**
large in scale
low quality
not realistic



## Multi-view images

**CO3D**
large in scale
No 3D GT



## Real-world 3D scans

**Google scanned objects**
high quality
real-world scans
household objects



## OmniObject3D

large-vocabulary
high quality
real-world scans

6K models from around 200 classes

Textured meshes

Point clouds    Rendered images

Real-captured videos

**Perception**    **Novel View Synthesis**    **Surface Reconstruction**    **Generation**

ModelNet → OmniObject3D → OmniObject3D-C

OOD Styles — OOD Corruptions

**Differences between CAD models and real-scanned objects**

**Common corruptions**

PointCloud-C (Ren et al. 2022)

Clean   Scale   Rotate   Jitter

Drop Global   Drop Local   Add Global   Add Local

ModelNet pretrained ⇨ OmniObject3D / OmniObject3D-C evaluation

Table 2. **Point cloud perception robustness analysis on Om-niObject3D with different architecture designs.** Models are trained on the ModelNet-40 dataset, with $OA_{Clean}$ to be their overall accuracy on the standard ModelNet-40 test set. $OA_{Style}$ on OmniObject3D evaluates the robustness to OOD styles. mCE on the corrupted OmniObject3D-C evaluates the robustness to OOD corruptions. Blue shadings indicate rankings. †: results on ModelNet-C [75]. Full results are presented in the supplementary materials.

| | $mCE^{\dagger} \downarrow$ | $OA_{Clean} \uparrow$ | $OA_{Style} \uparrow$ | $mCE \downarrow$ |
|---|---|---|---|---|
| DGCNN [92] | 1.000 | 0.926 | 0.448 | 1.000 |
| PointNet [71] | 1.422 | 0.907 | 0.466 | 0.969 |
| PointNet++ [72] | 1.072 | 0.930 | 0.407 | 1.066 |
| RSCNN [51] | 1.130 | 0.923 | 0.393 | 1.076 |
| SimpleView [30] | 1.047 | **0.939** | 0.476 | 0.990 |
| GDANet [99] | 0.892 | 0.934 | 0.497 | **0.920** |
| PAConv [98] | 1.104 | 0.936 | 0.403 | 1.073 |
| CurveNet [97] | 0.927 | 0.938 | **0.500** | 0.929 |
| PCT [32] | 0.925 | 0.930 | 0.459 | 0.940 |
| RPC [75] | **0.863** | 0.930 | 0.472 | 0.936 |

# Novel view synthesis (two settings)

☐ *Single-scene optimization models*

Multi-view images from **one scene** ⇒ **Train** ⇒ **The same scene Inference**



- NeRF (Mildenhall et al., 2021)
- Mip-NeRF (Barron et al., 2021)
- Plenoxels (Yu et al., 2021)

☐ *Generalizable models*

Multi-view images from scenes of **one category** or **across different categories**. ⇒ **Train** ⇒ **New scene with one or few views Inference**



- pixelNeRF (Yu et al., 2021)
- MVSNeRF (Chen et al., 2021)
- IBRNet (Wang et al., 2021)

# Surface reconstruction (two settings)

☐ *Multi-view image surface reconstruction*

# 3D object generation

*3D Object Generation*



*Interpolation across different categories*

# OmniObject3D: Text-to-3D Object



I want to generate a *toy dinosaur*.

I want to generate a *music box*.

I want to generate a *plaster statue*.

# Voxurf: Fast 3D Object Reconstruction

# Voxurf: Fast 3D Object Reconstruction

# Visual AIGC

# 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

LDR – Low Dynamic Range, [0, 255]

HDR – High Dynamic Range, [0, +∞]

# Text2Light: Text-to-3D Environment



"brown wooden dock on lake surrounded
by green trees during daytime"

**4K+ Resolution with High Dynamic Range**

"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

# Text2Light: Text-to-3D Environment

# SceneDreamer: Unbounded 3D Scene Generation



**In-the-wild 2D Image Collections**

**Photorealistic Unbounded 3D Scenes**

# SceneDreamer: Unbounded 3D Scene Generation



Multi-view consistent

Well-defined geometry

In-the-wild Image Collections

Photorealistic Unbounded 3D Scenes

Diverse scenes and styles

Infinite 3D World!

Generate with Different Styles

Free Camera Trajectories

# F2NeRF: Mobile 3D Scene Reconstruction

What if the input camera trajectory is very irregular? – We call that a "free" trajectory

# F2NeRF: Mobile 3D Scene Reconstruction

Adaptive warping method from input trajectories



(a) Space subdivision

(b) Local warping

(c) Different Hash function

Hash table

(d) Indexed with same Hash table

Indexed feature + View direction

Tiny-MLP

Density & Color

# F2NeRF: Mobile 3D Scene Reconstruction

# F2NeRF: Mobile 3D Scene Reconstruction

# Visual AIGC

Avatar

Object

Scene

**Foundation Model**

# Relighting4D: Relightable 3D Human



Prior works

Synthetic dataset

Light Stage data

Relighting4D uses **only** videos to relight dynamic human actors from free viewpoints
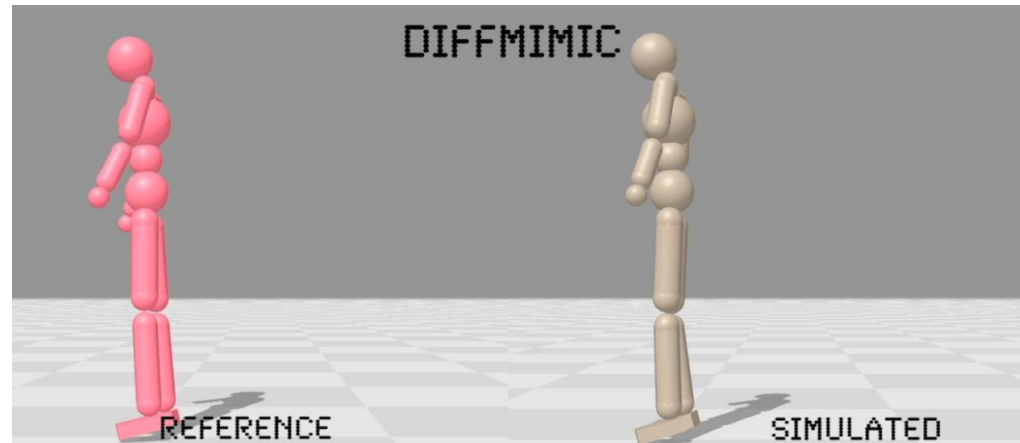
# Relighting4D: Relightable 3D Human



Video of human

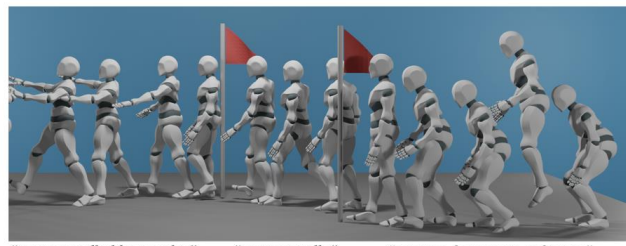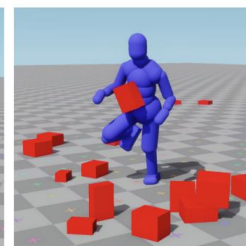Relight with different illuminations and free viewpoints

# Visual AIGC

Avatar

Foundation
Model

Object

Scene

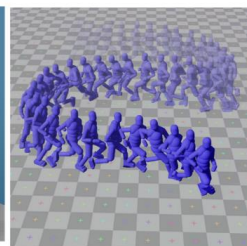# DiffMimic: Physically-Simulated Character

- Motion mimicking: let a **physically-simulated** character imitate a reference motion.
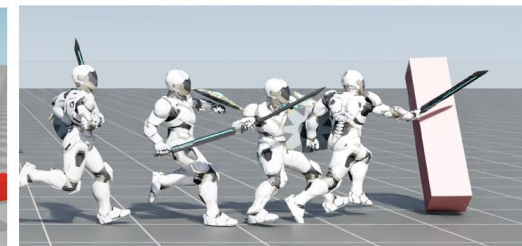


- A fundamental task for downstream animation applications.



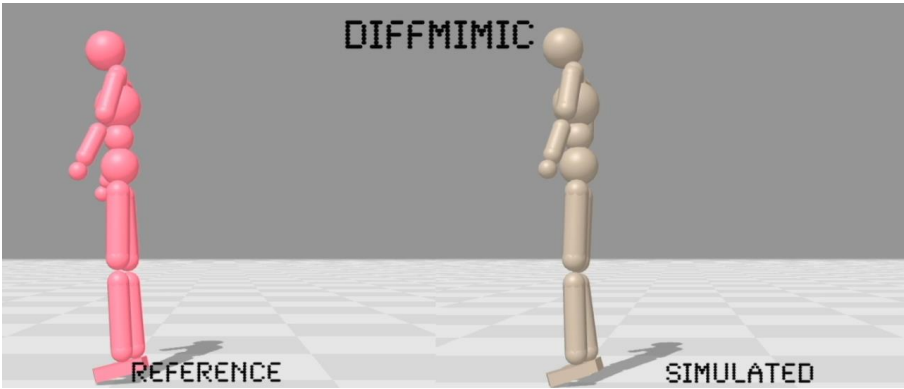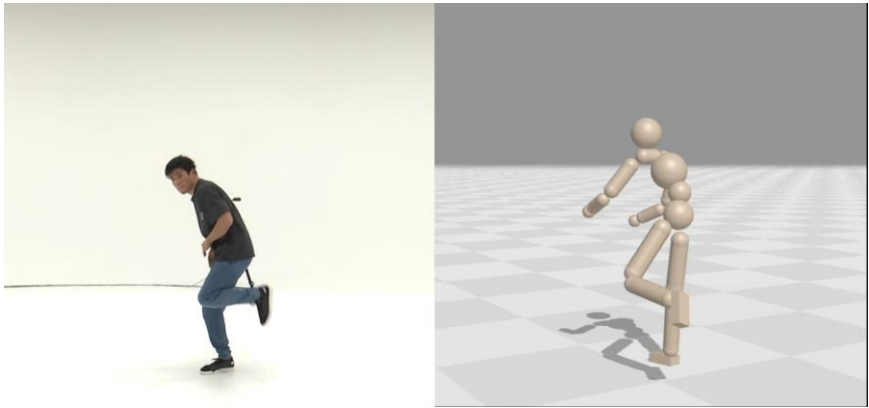Language-Conditioned Control        Responsive Control        Skill Composition

# DiffMimic: Physically-Simulated Character

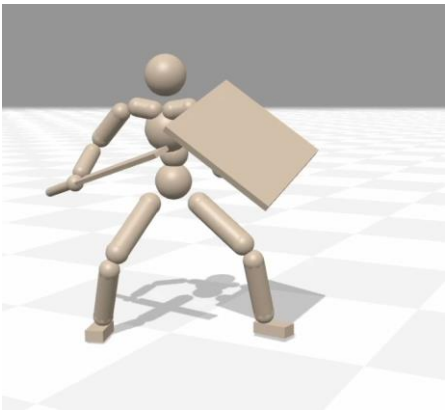| Motion | $T_{cycle}(s)$ | DeepMimic | Spacetime Bound | | Ours w/ RSI | |
|--------|-----------|-----------|-----------------|--|-------------|--|
| Back-Flip | 1.75 | 31.18 | 41.20 | +32.1% | 3.82 | -87.7% |
| Cartwheel | 2.72 | 30.45 | 17.35 | -43.0% | 4.72 | -84.5% |
| Walk | 1.25 | 23.80 | 4.08 | -79.5% | 1.55 | -93.5% |
| Run | 0.80 | 19.31 | 4.11 | -78.7% | 1.41 | -92.7% |
| Jump | 1.77 | 25.65 | 41.63 | +77.8% | 2.12 | -91.7% |
| Dance | 1.62 | 24.59 | 10.00 | -59.3% | 2.19 | -91.1% |

a) ~10x better sample efficiency compared to DeepMimic
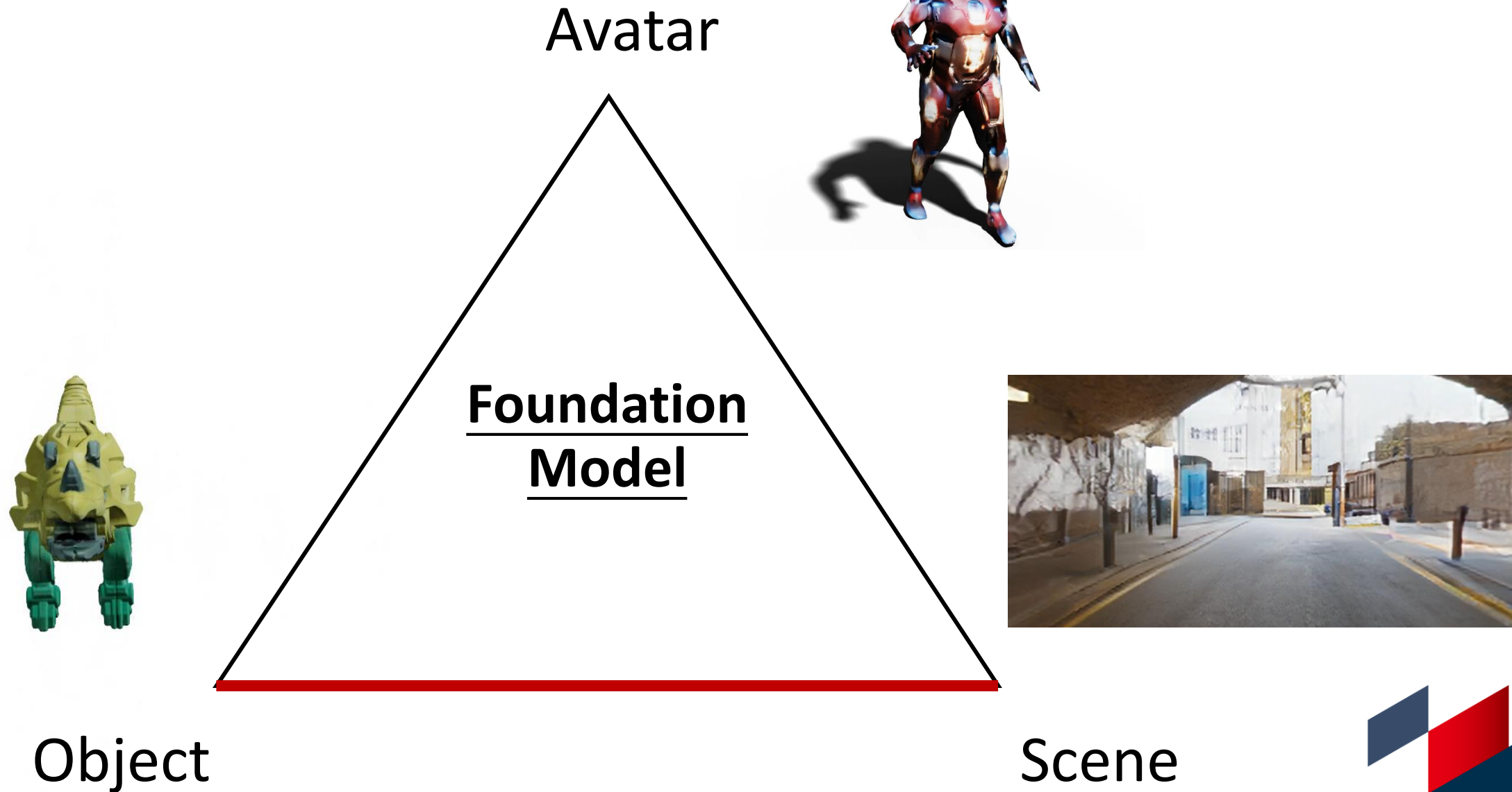


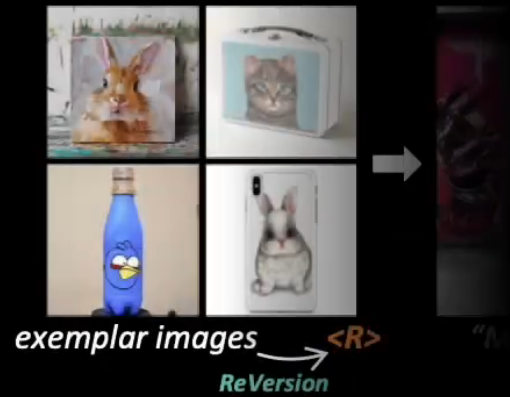b) Learning backflip in 5 minutes



c) Scalable



d) General



e) Robust

# Visual AIGC

# ReVersion: Object Relation Generation

# ReVersion: Object Relation Generation
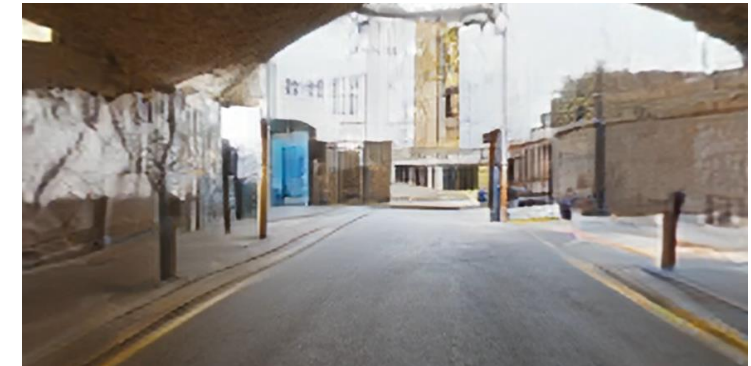
Visual Results: *ReVersion*

exemplar images → <R>
ReVersion

# Visual AIGC

Avatar



Thank You!

Object

Scene