# 3D Content Creation

## 3D content lays the foundation for broad applications


Avengers

Movies / VFX


GTA5

Gaming


Ready Player One

AR / VR / Metaverse



Manufacturing

## 3D assets: limited and expensive

# 3D Content Creation via Generative Models



$z \sim \mathcal{N}(0, I)$

Generative Model

Objverse, https://objaverse.allenai.org/

# Generative Models

**Diffusion Model**

Autoregressive Model

Implicit MLP
(Shap-E, 2023)

Vector set
(3DShape2VecSet, SIGGRAPH2023)

Feature maps
(LN3Diff, ECCV2024)

Encoder

Decoder

- Lack of explicit 3D-aware latent space for interactive editing.
- Lack of high-quality texture and efficient 3D VAE encoding from 2D inputs.

# GaussianAnything: Interactive Point Cloud Latent Diffusion for 3D Generation

Yushi Lan[1]    Shangchen Zhou[1]    Zhaoyang Lyu[2]    Fangzhou Hong[1]

Shuai Yang[3]    Bo Dai[2]    Xingang Pan[1]    Chen Change Loy[1]

[1]S-Lab, NTU Singapore    [2]Shanghai AI Lab    [3]Peking University

https://nirvanalan.github.io/projects/GA/

Lan et al, GaussianAnything: Interactive Point Cloud Latent Diffusion for 3D Generation, ICLR2025

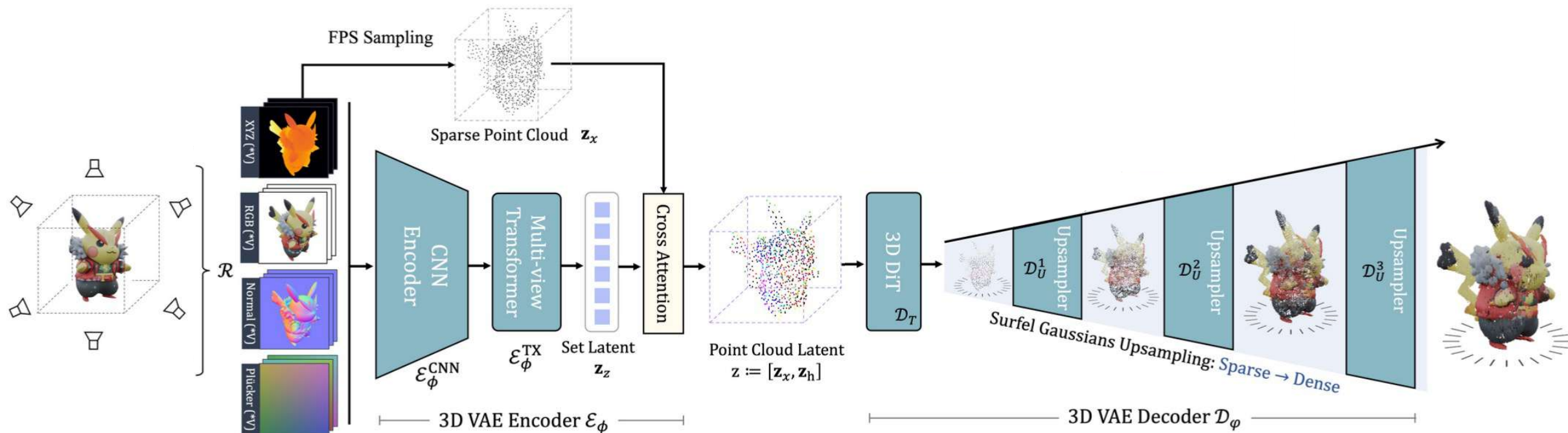# 3D VAE with Structured Latent

**Input**

Point cloud
+ multi-view RGB, Depth, Normal

**Latent**

Point cloud
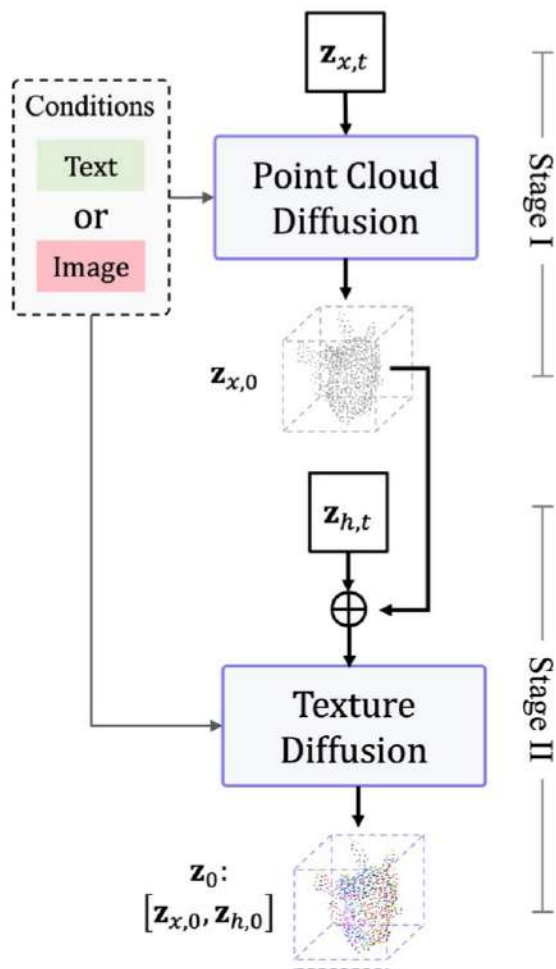+ Embedding for each point

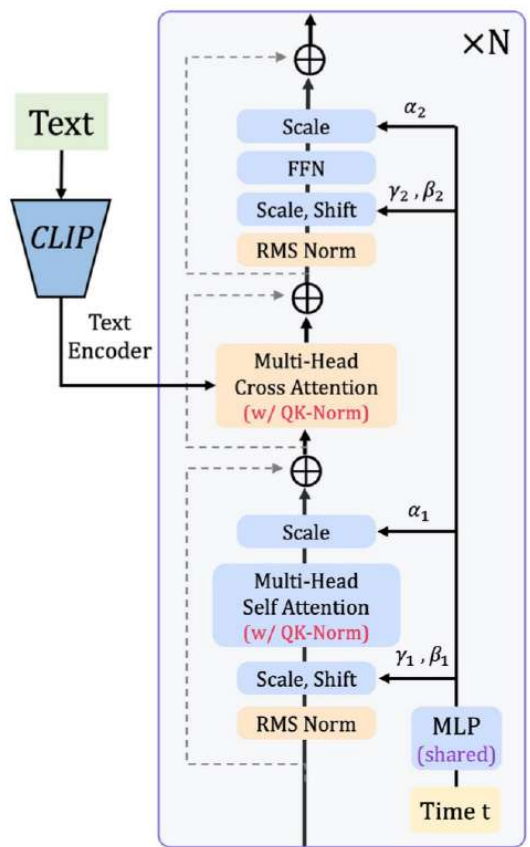**Output**

High-quality
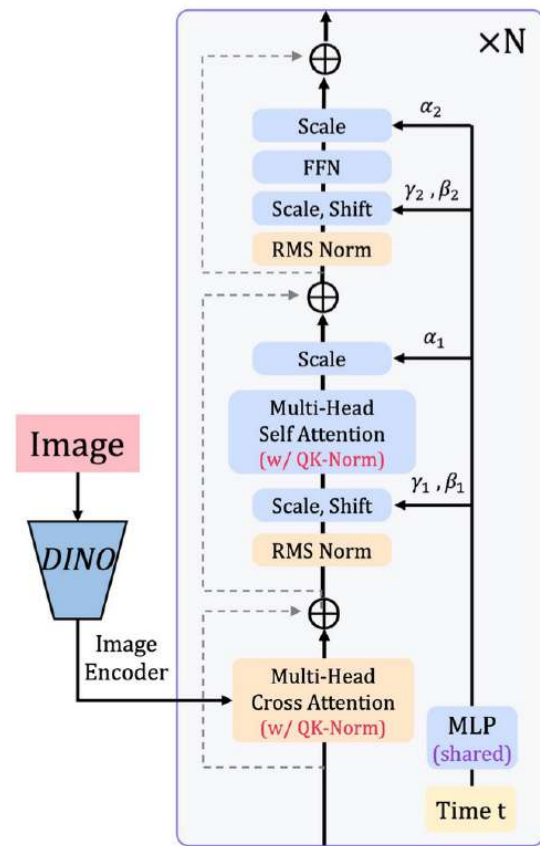surfel Gaussians



**Pipeline of the 3D VAE of GaussianAnything.**
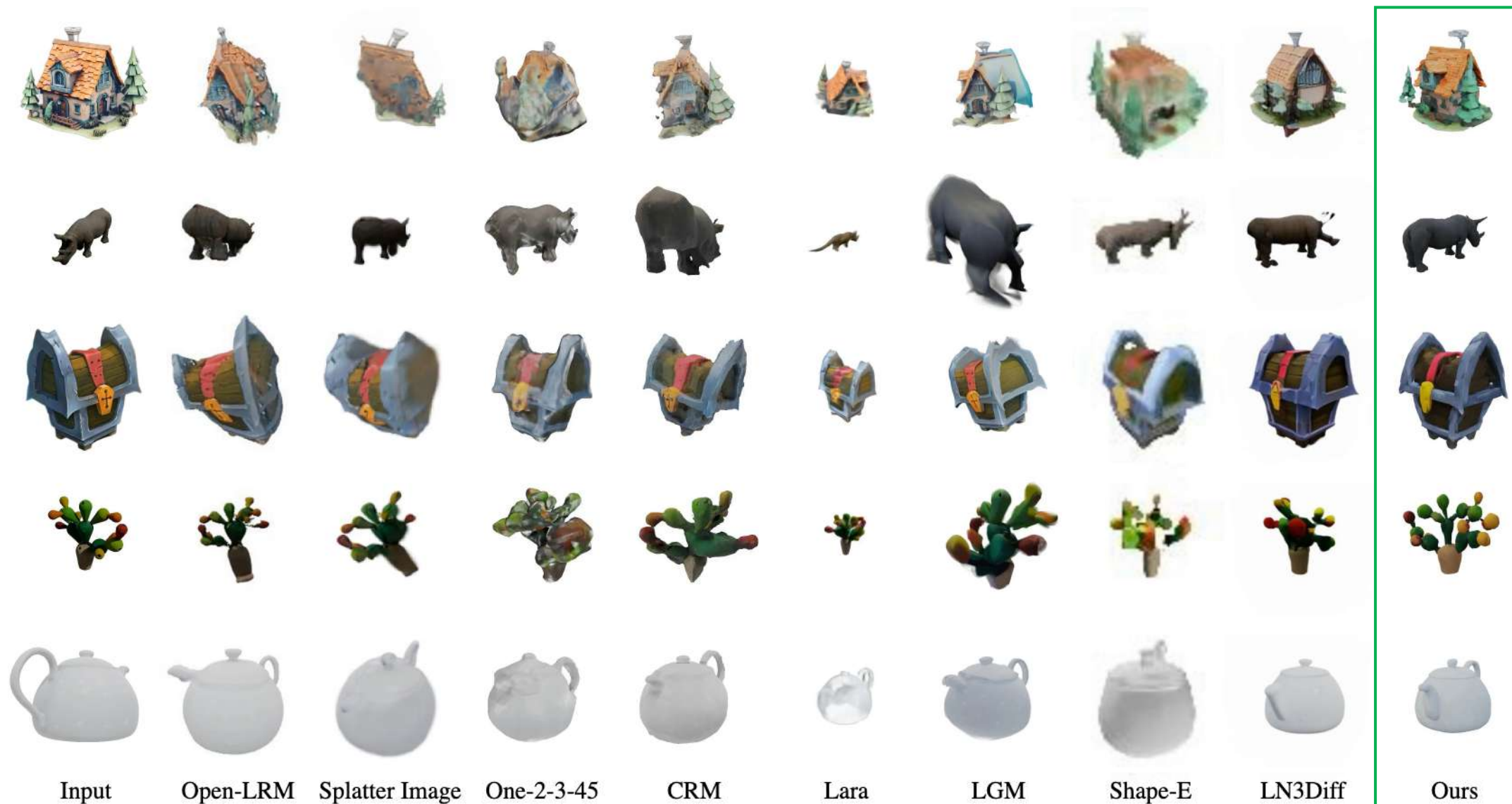
# Cascaded Native 3D Diffusion



Two-Stage Diffusion

(a) DiT Block (Text condition)

(b) DiT Block (Image condition)

Input · Open-LRM · Splatter Image · One-2-3-45 · CRM · Lara · LGM · Shape-E · LN3Diff · Ours

# Quantitative performance (Image-to-3D)

| Method | FID↓ | KID(%)↓ | MUSIQ↑ | P-FID↓ | P-KID(%)↓ | COV(%)↑ | MMD(‰)↓ |
|---|---|---|---|---|---|---|---|
| OpenLRM | 38.41 | 1.87 | 45.46 | 35.74 | 12.60 | 39.33 | 29.08 |
| Splatter-Image | 48.80 | 3.65 | 30.33 | 19.72 | 7.03 | 37.66 | 30.69 |
| One-2-3-45 (V=12) | 88.39 | 6.34 | 59.02 | 72.40 | 30.83 | 33.33 | 35.09 |
| CRM (V=6) | 45.53 | 1.93 | 64.10 | 35.21 | 13.19 | 38.83 | 28.91 |
| Lara (V=4) | 43.74 | 1.95 | 39.37 | 32.37 | 12.44 | 39.33 | 28.84 |
| LGM (V=4) | 19.93 | 0.55 | 54.78 | 40.17 | 19.45 | 50.83 | 22.06 |
| Shape-E | 138.53 | 11.95 | 31.51 | 20.98 | 7.41 | 61.33 | 19.17 |
| LN3Diff | 29.08 | 0.89 | 50.39 | 27.17 | 10.02 | 55.17 | 19.94 |
| **Ours** | 24.21 | 0.76 | 65.17 | 8.72 | 3.22 | 59.50 | 15.48 |

# Text-to-3D performance



A voxelized dog.

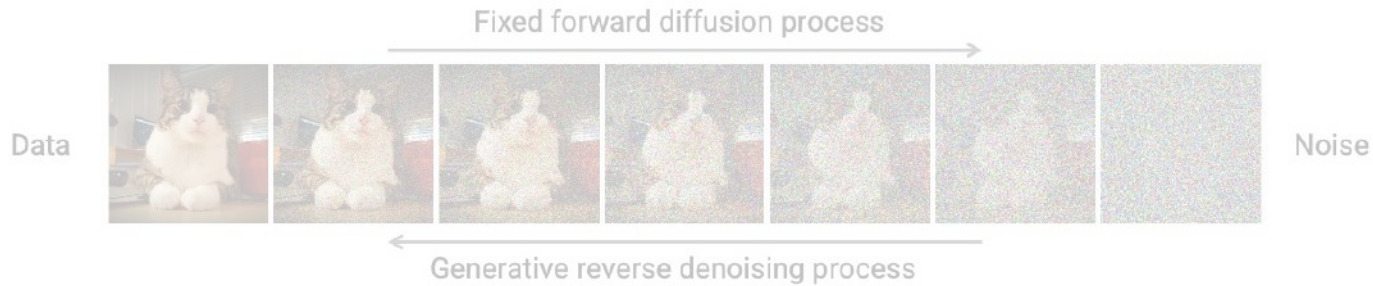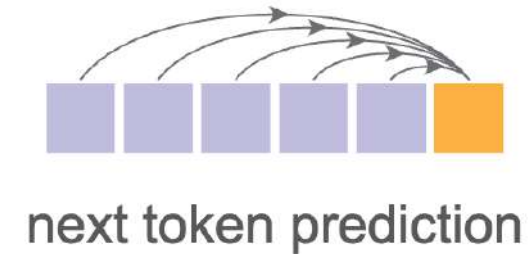An 18th century cannon.

Ours

3D Topia

LN3Diff

Shap-E

Point-E

# Interactive 3D Editing

# Generative Models

**Diffusion Model**

**Autoregressive Model**



Fixed forward diffusion process

Data

Noise

Generative reverse denoising process



next token prediction
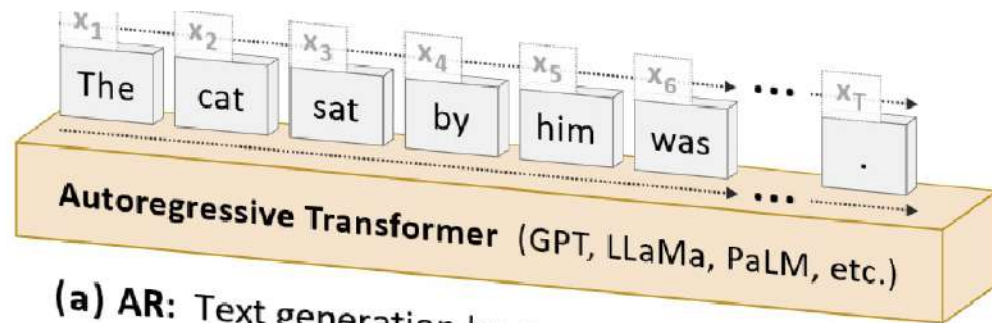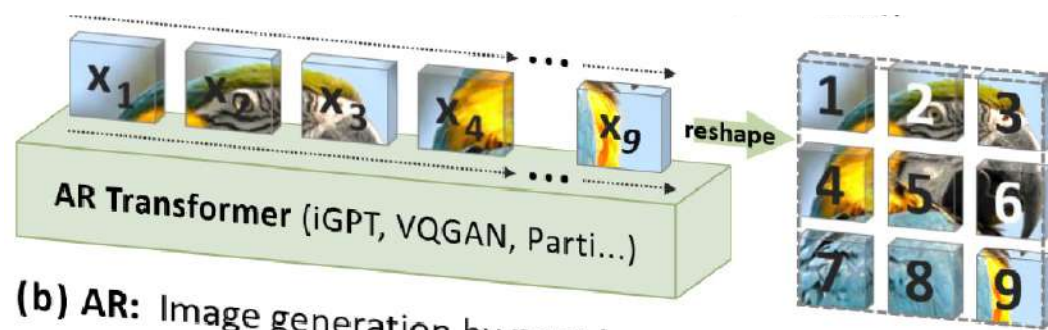
- A native 3D Diffusion Model, GaussianAnything:
  Structured latent space, better design, better performance

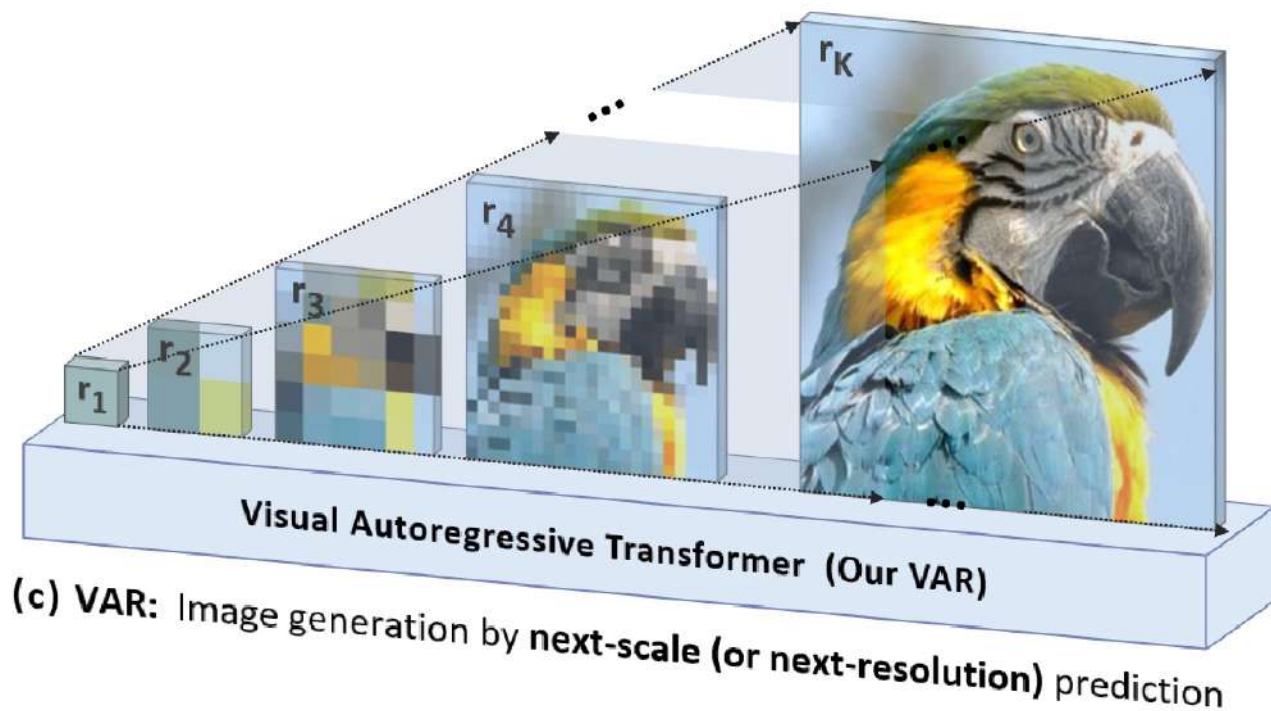Tian et al, Visual Autoregressive Modeling: Scalable Image Generation via Next-Scale Prediction, NeurIPS2024

Chen et al, SAR3D: Autoregressive 3D Object Generation and Understanding via Multi-scale 3D VQVAE, CVPR2025

# SAR3D -- Method



(a) 3D Generation

(b) 3D Understanding

Scale 2

Fast 3D generation (<1s)

Detailed 3D understanding

(a) Single Image to 3D

(b) Text to 3D

"a small brown cannon"

"A brown wooden table featuring visible legs and vintage look."

"A wooden chest"

"A coffee mug"

A sleek and aerodynamic blue and white racing car with a futuristic design, featuring racing stripes, a spoiler on the back, and a low profile.

A small, wooden house with a rectangular shape, staircase leading up to the entrance, and a patio area in front.

A wooden desk and chair set with a rectangular shape, featuring a simple and minimalistic design. The desk has a wooden top with a metal base, while the chair has a wooden seat and backrest.

A unique pair of black and green sunglasses with a slim and curved frame, featuring green lenses and a distinctive design.

**3D Captioning**. Given a 3D model, our method can generate captions that contain both category and details.

# SAR3D -- Comparison

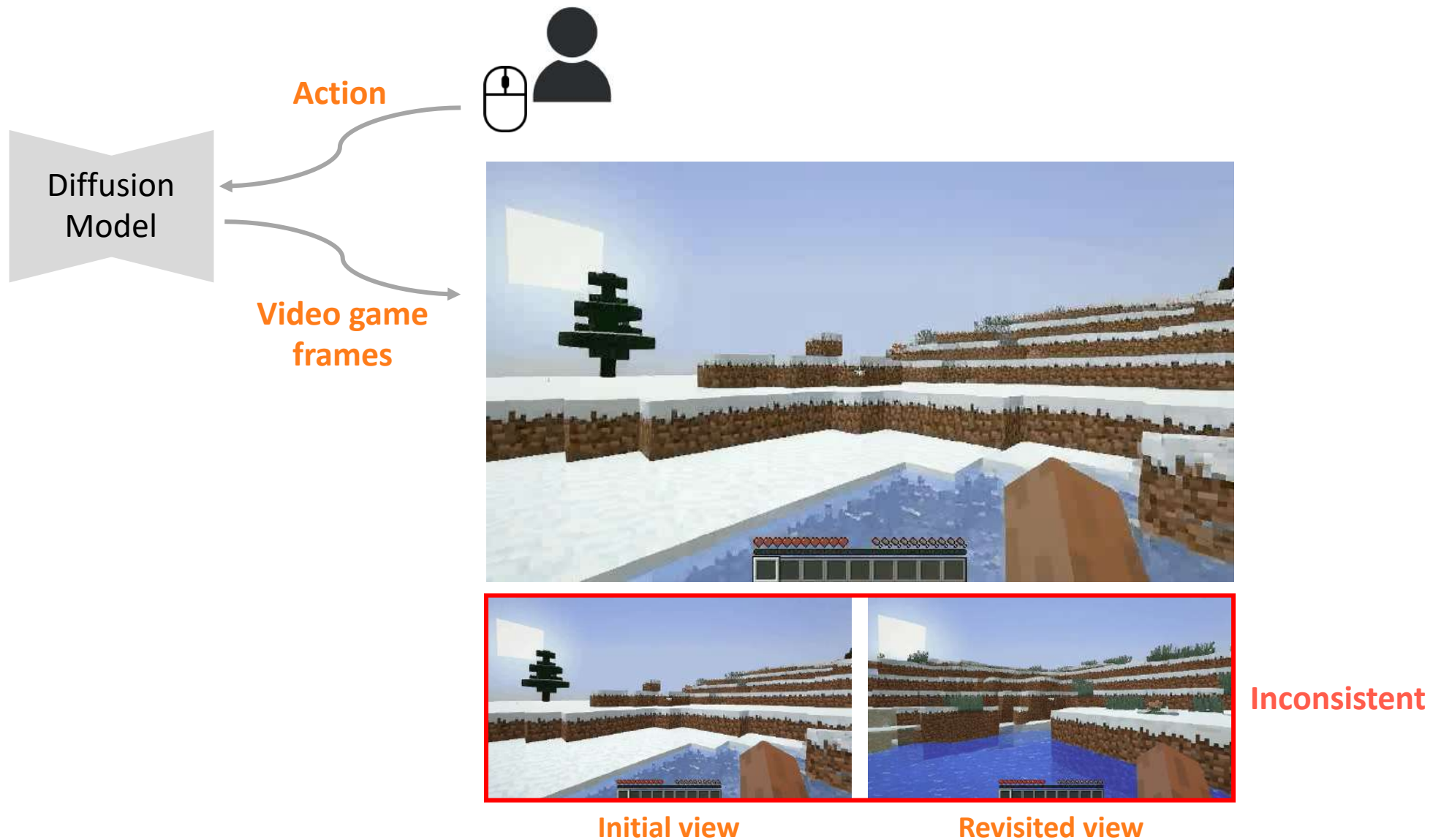| Method | FID↓ | KID(%)↓ | MUSIQ↑ | COV(%)↑ | MMD(‰)↓ | Latency-V100 (s) ↓ |
|---|---|---|---|---|---|---|
| Splatter-Image | 48.80 | 3.65 | 30.33 | 37.66 | 30.69 | 0.83 |
| OpenLRM | 38.41 | 1.87 | 45.46 | 39.33 | 29.08 | 7.21 |
| One-2-3-45 (V=12) | 88.39 | 6.34 | 59.02 | 33.33 | 35.09 | 59.23 |
| Lara (V=4) | 43.74 | 1.95 | 39.37 | 39.33 | 28.84 | 11.93 |
| CRM (V=6) | 45.53 | 1.93 | 64.10 | 38.83 | 28.91 | 22.10 |
| LGM (V=4) | 19.93 | 0.55 | 54.78 | 50.83 | 22.06 | 3.87 |
| Shap-E | 138.53 | 11.95 | 31.51 | 61.33 | 19.17 | 9.54 |
| LN3Diff | 29.08 | 0.89 | 50.39 | 55.17 | 19.94 | 7.51 |
| GaussianAnything | 24.21 | 0.76 | 65.17 | 59.50 | 15.48 | 15.02 |
| **SAR3D**-NeRF | 22.55 | 0.42 | 65.77 | 74.17 | 13.63 | 1.64 |
| **SAR3D**-Flexicubes | 27.30 | 0.63 | 67.24 | 71.50 | 15.25 | 2.92 |

- Autoregressive Model for 3D generation can perform as well as diffusion models while being more efficient

# Generative AI as 3D Game Engine?



Minecraft powered purely by Generative AI
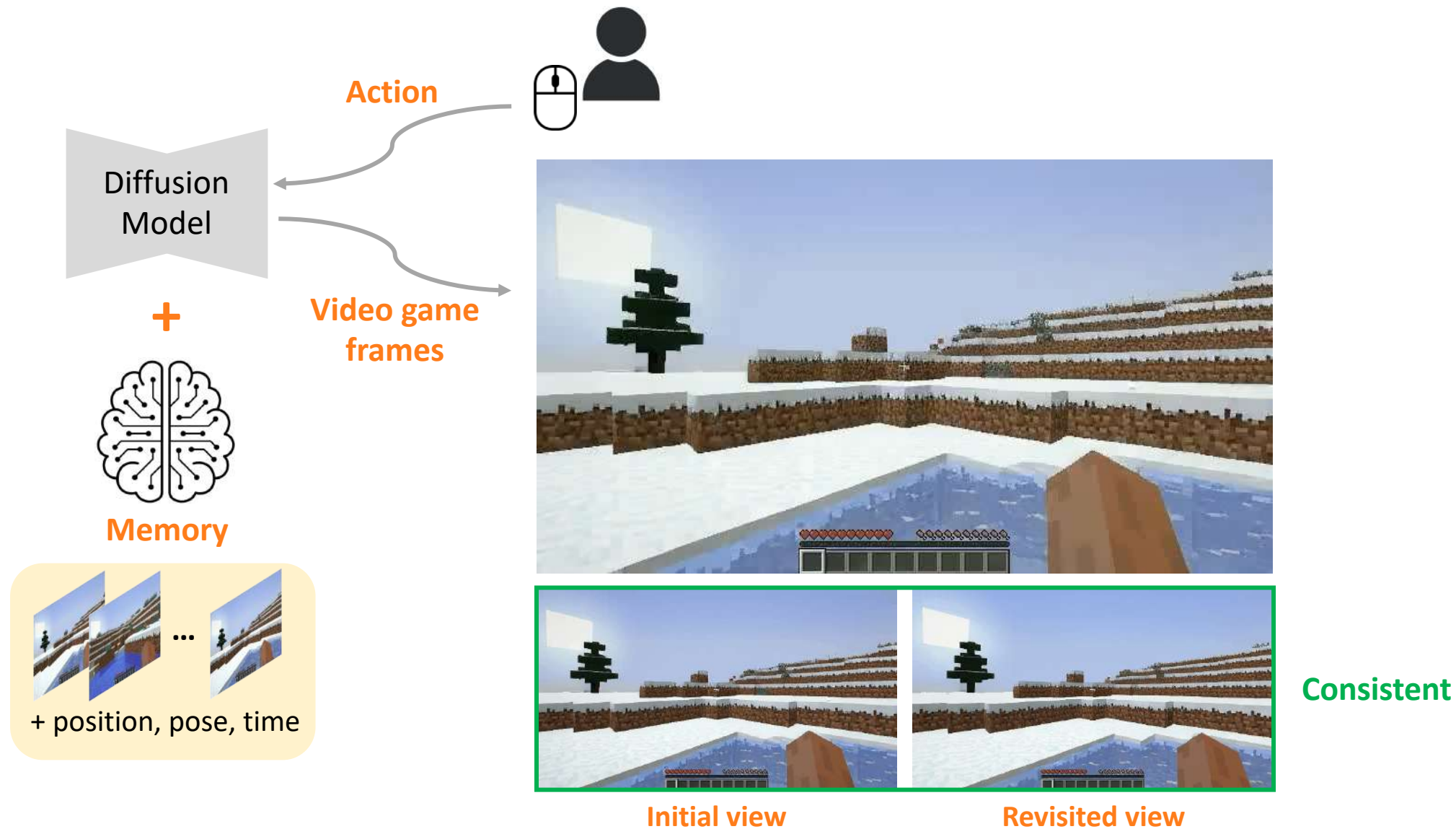
Xiao et al, To appear on arXiv, 2025

# Generative AI as 3D Game Engine



Action

Diffusion Model

Video game frames

Inconsistent

Initial view

Revisited view

Oasis: A Universe in a Transformer https://oasis-model.github.io/

# Generative AI as 3D Game Engine



Xiao et al, To appear on arXiv, 2025

# World Generation with Memory



History frames
(long-term memory)

Short past
(within temporal attention window)

Next frame
to be predicted

**+ position, pose, time**

Diffusion
Model

Most relevant
memories

**Retrieve**
(FOV overlap, time)

**+ position,
pose, time**

**Memory attention**
(pose+time aware)

FOV of View 1
FOV of View 2
Overlap

View 1    View 2

# Generative AI as 3D Game Engine



Action

Diffusion Model

+

Memory

Video game frames

Initial view

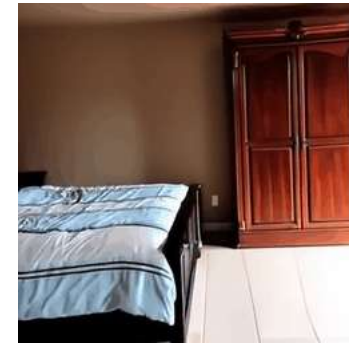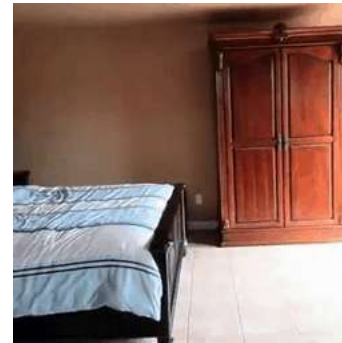Revisited view

# Real Scene Results



Initial view     Revisited view

**w/o Memory**

Initial view     Revisited view

**w/ Memory**

# Real Scene Results



Initial view     Revisited view

Initial view     Revisited view

# Conclusion



Scale 2

Fast 3D generation (<1s)

- A native 3D Diffusion Model, GaussianAnything: structured latent space, better design, better performance
- Autoregressive Model for 3D generation can perform as well as diffusion models while being more efficient
- When building 3D playable worlds via video diffusion models, **Memory** is important!

## Open problems

- 3D object -> Rigging -> **Animation**
- 3D **scene** generation
- **CAD** generation, 3D object to CAD
- **Physics**-aware generation

# Thank you!