

Large Reconstruction Model (LRM)



Introduction

Task: X to 3D

• Image-to-3D

Text-to-3D



• Multi-view-to-3D, Video-to-3D, ...

Results from [Xu & Shi et al. '24]

Reconstruction VS

Generation

What is a Large Reconstruction Model?



Feed-Forward Network

trained on large data to generalize



Input: single image, text, ...

Output: Radiance field representation, e.g., NeRF, 3DGS



Large Reconstruction Model (LRM)

PixelNeRF



PixelNeRF















7

Input image PixelNeRF

Prof. Niessner

[Yu et al. '21] pixelNeRF: Neural Radiance Fields from One or Few Images





1. Extract features from input image with CNN encoder

Input Image

2. Shoot camera rays from pixels of the target view

Target Image

1. Extract features from input image with CNN encoder

Image Encoder

Prof. Niessner

Input Image

2. Shoot camera rays from pixels of the target view

Target Image

1. Extract features from input image with CNN encoder

Image Encoder

3. Determine features for sample points along the rays:

- Project 3D points to image plane.
- Bilinearly interpolate image features





→ Volume rendering and supervision as in NeRF.

PixelNeRF – Multi-View Input





PixelNeRF – Architecture



LRM: Large Reconstruction Model



Prof. Niessner

LRM – Image Encoder Patchify image into sequence Pretrained Single input image Image encoder Image features Dim: 512 x 512 x 3 12 Layers, Dim: 768, ViT (DINO) Dim: (32 x 32) x 768 MLP Self Conv Patch-wise Self-attention feature tokens between patches

LRM: Large Reconstruction Model



Prof. Niessner



20

LRM – Image-to-Triplane Decoder



2 different conditional operations

 $\gamma, \beta = \mathrm{MLP}^{\mathrm{mod}}(\tilde{c})$ ModLN_c $(f_j) = \mathrm{LN}(f_j) \cdot (1 + \gamma) + \beta$

Modulation with camera features
→ control orientation and
distortion of the whole shape

Cross-attention with input image patches
 → control fine-grained geometric and color information

[Hong et al. '24] LRM: Large Reconstruction Model for Single Image to 3D 21

LRM: Large Reconstruction Model



Prof. Niessner

[Hong et al. '24] LRM: Large Reconstruction Model for Single Image to 3D 22

LRM – Triplane NeRF



LRM – Results



Prof. Niessner

[Hong et al. '24] LRM: Large Reconstruction Model for Single Image to 3D 24

GRM: Large Gaussian Reconstruction Model

Single image to 3D

Text to 3D



Multi-view to 3D





Prof. Niessner

[Xu & Shi et al. '24] GRM: Large Gaussian Reconstruction Model for Efficient 3D Reconstruction and Generation

GRM: Large Gaussian Reconstruction Model



Prof. Niessner

[Xu & Shi et al. '24] GRM: Large Gaussian Reconstruction Model for Efficient 3D Reconstruction and Generation 26

GRM: ViT Encoder



1. Concatenate image with camera poses as Plücker rays 2. Encode with CNN tokenizer to create one sequence from all 4 images $4 \times \frac{H}{16} \times \frac{W}{16}$ 3. Append learnable image position encodings. 4. Series of self-attention layers attending to all tokens across the input views.

GRM: Large Gaussian Reconstruction Model



Prof. Niessner

[Xu & Shi et al. '24] GRM: Large Gaussian Reconstruction Model for Efficient 3D Reconstruction and Generation 28

GRM: Transformer-based Upsampler

- Multiple upsample blocks progressively upsample by factor 2 until reaching $H \times W$
 - Quadruple number of channels
 - Double spatial dimension with PixelShuffle



GRM: Transformer-based Upsampler

- Multiple upsample blocks progressively upsample by factor 2 until reaching $H \times W$
 - Quadruple number of channels
 - Double spatial dimension with **PixelShuffle**



https://nico-curti.github.io/NumPyNet/NumPyNet/layers/pixelshuffle_layer.html

GRM: Transformer-based Upsampler

- Windowed Self-Attention: balance between need for non-local multi-view information aggregation and feasible computation cost
- Separate linear heads produce Gaussian features



GRM: Large Gaussian Reconstruction Model



Prof. Niessner

[Xu & Shi et al. '24] GRM: Large Gaussian Reconstruction Model for Efficient 3D Reconstruction and Generation 32

GRM: Pixel-aligned 3D Gaussians

- Predict one Gaussian per pixel, i.e., one Gaussian attribute map $H \times W \times C$ per input image, C = 12, depth (1), rotation (4), Gaussians scale (3), opacity (1), rgb (3).
- \rightarrow Pixel-aligned Gaussians are easier to learn
- → 3D Gaussians render in real-time

Image from [Charatan et al. '24]

GRM: Pixel-aligned 3D Gaussians

- Backproject pixel-aligned Gaussians to single 3DGS representation
- 2. Render as in 3DGS
- Supervise novel views with image (color & perceptual) and mask loss



Prof. Niessner

GRM - Results



35

Long-LRM: Long-sequence Large Reconstruction Model



→ from 32 images in 1.3s

Long-LRM: Long-sequence Large **Reconstruction Model**



Wide-coverage Gaussian Reconstruction

Prot. Niessner

Long-LRM – Hybrid Block

- Combine Mamba2 blocks & transformer blocks →
 better scalability to higher resolution & denser views
- Hybrid block = 7 Mamba2 blocks + 1 transformer block
 - L = Sequence Length



Due to global self-attention

 $\bigcirc (| \land 2)$

Long-LRM – Mambaz Block

- Mamba2 block [Dao & Gu ,24] is designed for language tasks → scans through sequence in one direction → suboptimal for images.
- Vision Mamba [Zhu et al. '24] take bi-directional scans over the concatenated token sequence



Prof. Niessner

Long-LRM: Long-sequence Large **Reconstruction Model**



Wide-coverage Gaussian Reconstruction

Prot. Niessner

Long-LRM – Token Merge

- 1. Reshape input sequence $L \times D$ back to $N \times \frac{H}{p} \times \frac{w}{p} \times D$ where N = #images, p = patch size
- 2. Apply 2D convolution, kernel size 2, stride 2, resulting in $N \times \frac{H}{2p} \times \frac{W}{2p} \times D'$
- 3. Resape back to $\frac{L}{4} \times D'$
- \rightarrow Reduces sequence length to 1/4



Prof. Niessner

Long-LRM: Long-sequence Large **Reconstruction Model**



Wide-coverage Gaussian Reconstruction

Prot. Niessner

Long-LRM – Decode to Gaussians

Decode output tokens to per-pixel Gaussian
 parameters

Novel View Synthesis

 At training-time prune to fixed number of Gaussians, at test-time prune by opacity → improve efficiency at high resolution and increased views



Wide-coverage Gaussian Reconstruction

Prof. Niessner

$$\begin{aligned} & \text{Long-LRM} - \text{Training Objective} \\ \mathcal{L} &= \mathcal{L}_{\text{image}} + \lambda_{\text{opacity}} \cdot \mathcal{L}_{\text{opacity}} + \lambda_{\text{depth}} \cdot \mathcal{L}_{\text{depth}} \\ \mathcal{L}_{\text{image}} &= \frac{1}{M} \sum_{i=1}^{M} \left(\text{MSE} \left(\mathbf{I}_{i}^{\text{gt}}, \mathbf{I}_{i}^{\text{pred}} \right) + \lambda \cdot \text{Perceptual} \left(\mathbf{I}_{i}^{\text{gt}}, \mathbf{I}_{i}^{\text{pred}} \right) \right) \\ \mathcal{L}_{\text{depth}} &= \frac{1}{M} \sum_{i=1}^{M} \text{Smooth-L1} \left(\mathbf{D}_{i}^{\text{da}}, \mathbf{D}_{i}^{\text{pred}} \right) & \text{Depth regularizer} \\ \mathcal{L}_{\text{opacity}} &= \frac{1}{N} \sum_{i=1}^{N} |o_{i}| & \text{Opacity regularizer to reduce the} \\ \end{aligned}$$

Prof. Niessner

Long-LRM – Results







Reconstruction Model for Wide-coverage Gaussian Splats

Avat3r: Generalizable 3D Head Avatars



arXiv'25 [Kirschstein et al.] Avat3r

Avat3r: Generalizable 3D Head Avatars



arXiv'25 [Kirschstein et al.] Avat3r







Large Animatable Reconstruction Model for High-fidelity 3D Head Avatars

Tobias Kirschstein^{1,2} - Javier Romero² - Artem Sevastopolsky^{1,2} - Matthias Nießner¹ - Shunsuke Saito²

¹Technical University of Munich ²M

²Meta Reality Labs













arXiv'25 [Kirschstein et al.] Avat3r













Zero-shot 3D facial animation

Reading Homework

- [Hong et al. '24] LRM: Large Reconstruction Model for Single Image to 3D
 - <u>https://arxiv.org/pdf/2311.04400</u>

Literature

- [Yu et al. '21] pixelNeRF: Neural Radiance Fields from One or Few Images
- [Hong et al. '24] LRM: Large Reconstruction Model for Single Image to 3D
- [Xu & Shi et al. '24] GRM: Large Gaussian Reconstruction Model for Efficient 3D Reconstruction and Generation

Literature

- [Chen et al. 2024] Long-LRM: Long-sequence Large Reconstruction Model for Wide-coverage Gaussian Splats
- arXiv '25 [Kirschstein et al.] Avat3r



Thanks for watching!