

LARM: A Large Articulated-Object Reconstruction Model

SYLVIA YUAN, University of California San Diego, USA

RUOXI SHI, University of California San Diego, USA and Hillbot Inc., USA

XINYUE WEI, University of California San Diego, USA and Hillbot Inc., USA

XIAOSHUAI ZHANG, Hillbot Inc., USA

HAO SU, University of California San Diego, USA and Hillbot Inc., USA

MINGHUA LIU, Hillbot Inc., USA

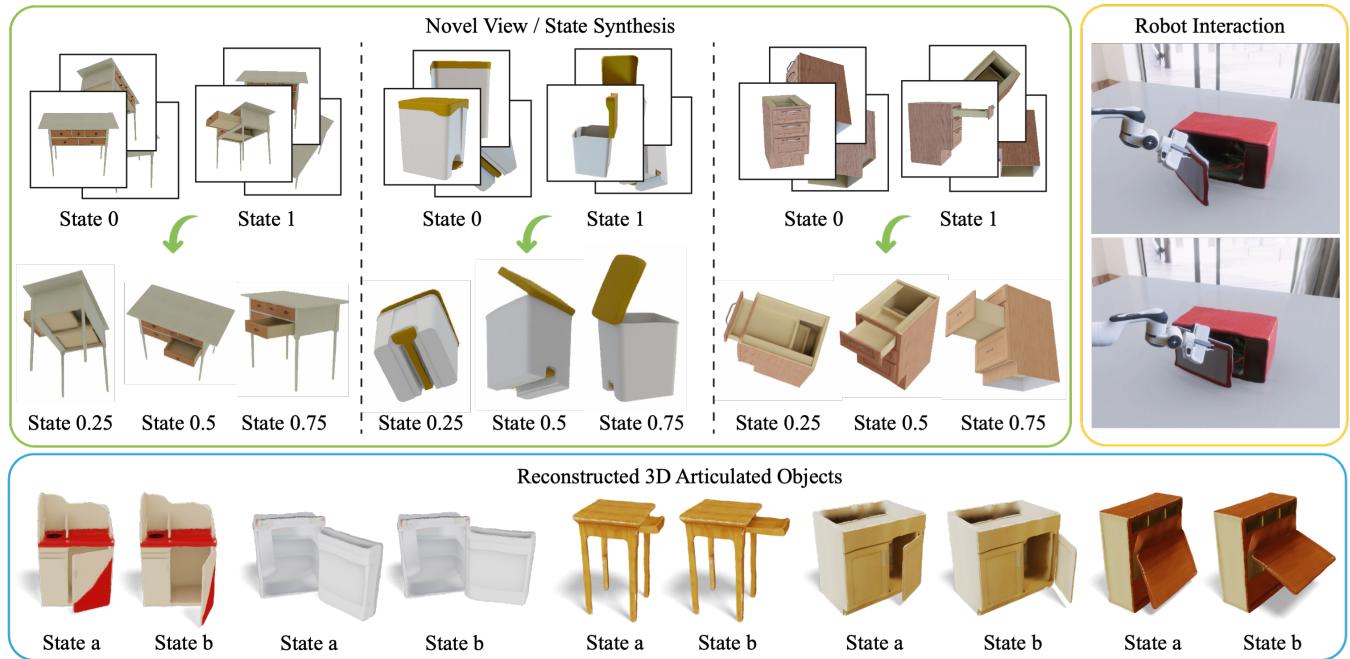


Fig. 1. We introduce LARM, a feedforward model for 3D articulated object reconstruction. Given sparse-view images of an articulated object in two different joint states (e.g., drawer open and closed), LARM can synthesize views from arbitrary camera poses and novel joint configurations (e.g., the drawer half-open), enabling efficient generation of continuous articulation and viewpoint variations. Beyond novel view and state synthesis, LARM also supports explicit 3D mesh reconstruction. Unlike many recent methods that rely on part retrieval or template meshes—approaches that often fail to capture fine-grained geometry and lack support for realistic textures—LARM reconstructs high-quality, textured, articulated 3D objects that closely match the input images in terms of geometric detail, visual fidelity, and articulation accuracy. LARM enables a wide range of downstream applications, including high-fidelity digital twin reconstruction from casual real-world captures, which is crucial for large-scale robotics learning in simulation.

Authors' Contact Information: Sylvia Yuan, University of California San Diego, USA, liyuan@ucsd.edu; Ruoxi Shi, University of California San Diego, USA and Hillbot Inc., USA, r8shi@ucsd.edu; Xinyue Wei, University of California San Diego, USA and Hillbot Inc., USA, xiwei@ucsd.edu; Xiaoshuai Zhang, Hillbot Inc., USA, x@hillbot.ai; Hao Su, University of California San Diego, USA and Hillbot Inc., USA, haosu@ucsd.edu; Minghua Liu, Hillbot Inc., USA, m@hillbot.ai.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

SA Conference Papers '25, Hong Kong, Hong Kong

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 979-8-4007-2137-3/2025/12

<https://doi.org/10.1145/3757377.3763844>

Modeling 3D articulated objects with realistic geometry, textures, and kinematics is essential for a wide range of applications. However, existing optimization-based reconstruction methods often require dense multi-view inputs and expensive per-instance optimization, limiting their scalability. Recent feedforward approaches offer faster alternatives but frequently produce coarse geometry, lack texture reconstruction, and rely on brittle, complex multi-stage pipelines. We introduce LARM, a unified feedforward framework that reconstructs 3D articulated objects from sparse-view images by jointly recovering detailed geometry, realistic textures, and accurate joint structures. LARM extends LVSM—a recent novel view synthesis (NVS) approach for static 3D objects—into the articulated setting by jointly reasoning over camera pose and articulation variation using a transformer-based architecture, enabling scalable and accurate novel view synthesis. In addition, LARM generates auxiliary outputs such as depth maps and part masks to facilitate explicit 3D mesh extraction and joint estimation. Our pipeline eliminates the need for dense supervision and supports high-fidelity reconstruction

across diverse object categories. Extensive experiments demonstrate that LARM outperforms state-of-the-art methods in both novel view and state synthesis as well as 3D articulated object reconstruction, generating high-quality meshes that closely adhere to the input images. Our project page is at <https://sylviayuan-sy.github.io/larm-site/>.

CCS Concepts: • Computer vision problems → Reconstruction; • Image and video acquisition → 3D imaging; • Computer vision tasks → Vision for robotics; • Animation → Motion processing.

ACM Reference Format:

Sylvia Yuan, Ruoxi Shi, Xinyue Wei, Xiaoshuai Zhang, Hao Su, and Minghua Liu. 2025. LARM: A Large Articulated-Object Reconstruction Model . In *SIGGRAPH Asia 2025 Conference Papers (SA Conference Papers '25), December 15–18, 2025, Hong Kong, Hong Kong*. ACM, New York, NY, USA, 17 pages. <https://doi.org/10.1145/3757377.3763844>

1 Introduction

3D articulated objects—from household appliances and furniture to complex mechanical assemblies—pervade everyday environments, yet faithfully modeling their geometry and kinematic structures remains a labor-intensive endeavor. Traditional pipelines rely on skilled artists who manually build and rig each asset, incurring substantial time and monetary costs, and large-scale datasets with detailed articulation labels are correspondingly scarce. Accurate reconstruction of such objects is, however, a prerequisite for high-fidelity robotic simulation and digital-twin systems, AR/VR and animation, where realistic geometry, appearance, and kinematics enable robust robotics training and low sim-to-real gap deployment.

Prior work has made significant strides in analyzing articulated objects [Liu et al. 2024b], particularly through tasks like joint axis prediction and joint state estimation [Fu et al. 2024; Liu et al. 2023a]. Beyond analysis, several optimization-based reconstruction methods—often built on NeRFs or 3D Gaussian Splatting—have been developed to recover articulated geometry [Liu et al. 2023c; Tseng et al. 2022; Wei et al. 2022]. While effective, these approaches are computationally expensive, require dense multi-view inputs, and rely on per-instance optimization with limited shape-level priors, hindering scalability and generalization. More recently, feedforward models have emerged for articulated object reconstruction or generation from various inputs [Chen et al. 2024; Le et al. 2024; Liu et al. 2024a]. These input conditions include object categories, articulation graphs, single images, and even text. However, many of these models [Chen et al. 2024; Le et al. 2024; Liu et al. 2024a] oversimplify geometry with bounding boxes, template meshes, or retrieved parts from a small database, leading to coarse and unrealistic results. Many also lack texture reconstruction, further reducing visual fidelity. Consequently, they struggle to produce meshes that accurately reflect user inputs, limiting their applicability to use cases like digital twins and high-fidelity simulation. A few recent feedforward methods [Kawana and Harada 2023; Zhang et al. 2021] model geometry with implicit fields (e.g., SDFs), but most of these methods involve complex, multi-stage pipelines, and still lack high-quality texture modeling.

In contrast, we propose LARM, a feedforward framework for reconstructing 3D articulated objects. Unlike prior methods, LARM jointly recovers detailed geometry, realistic textures, and accurate articulated structure from sparse-view inputs, eliminating the need

for dense-view observations. Built on a simple and unified architecture, it is fast, scalable, and capable of learning rich shape priors across a wide variety of articulated objects.

Our approach draws inspiration from LVSM [Jin et al. 2024], a recent method that demonstrates a transformer model with minimal 3D inductive bias can synthesize novel views of static objects. We extend this idea to articulated settings by enabling the model to reason over both camera pose and articulation variation. As shown in Figure 1, LARM takes as input sparse-view posed images taken from two articulation states (e.g., a rest state and a maximum state), each concatenated with its corresponding joint state. These inputs are processed by a transformer alongside target-view Plücker rays and a target joint state to predict the image tokens of the desired view. In doing so, LARM is trained to synthesize novel views conditioned on arbitrary camera poses and joint configurations, effectively capturing both viewpoint and articulation changes. While novel-view synthesis alone already supports a range of downstream applications, we further extend LARM to output auxiliary signals—including depth maps, foreground masks, and part segmentation masks—to facilitate explicit 3D reconstruction. Finally, we introduce post-processing modules that leverage these outputs for articulation joint estimation and explicit mesh extraction, enabling the full reconstruction of textured, articulated 3D assets.

We compare LARM with recent state-of-the-art methods for 3D articulated object reconstruction and generation, including both optimization-based and feedforward approaches. We evaluate performance on both novel view and state synthesis as well as 3D articulated object reconstruction tasks, where LARM significantly outperforms baseline methods across all tasks and metrics. Compared to baselines—especially retrieval-based ones—LARM demonstrates much better adherence to the input images in terms of both geometric and texture details. We also test our pipeline using casually captured images from handheld devices, showcasing the model’s potential for real-world applications such as AR/VR and large-scale digital twin generation for robot learning.

2 Related Work

2.1 Understanding of 3D Articulated Objects

Understanding the structure and motion of 3D articulated objects has been a central topic in both computer vision and robotics. Many efforts have been devoted to identifying movable parts from a single image [Jiang et al. 2022b; Sun et al. 2023] or from video sequences [Qian et al. 2022]. A large body of literature focuses on estimating joint parameters—such as rotation axes, pivot locations, and articulation angles—across different input modalities, including point clouds [Fu et al. 2024; Jiang et al. 2022a; Li et al. 2020; Liu et al. 2023a; Wang et al. 2019; Xu et al. 2022; Yan et al. 2020], RGB-D data [Abbatematteo et al. 2019; Che et al. 2024; Hu et al. 2017; Jain et al. 2022, 2021; Liu et al. 2022], videos [Liu et al. 2020], and 4D dynamic point clouds [Liu et al. 2023b]. In addition to static scene understanding, active perception techniques have been introduced to enhance the precision of articulation estimation [Yan et al. 2023; Zeng et al. 2024]. Furthermore, several works explore the temporal dynamics of articulated objects, addressing challenges such as continuous pose tracking and motion estimation [Heppert et al.

2022; Weng et al. 2021]. Recently, vision-language models have been fine-tuned to support joint estimation tasks [Huang et al. 2024].

2.2 Optimization-Based Articulated Object Reconstruction

In addition to detecting and understanding 3D articulated objects, a substantial body of work focuses on reconstructing them from multi-view images. These methods typically follow a per-instance optimization paradigm, extending frameworks such as Neural Radiance Fields (NeRF) [Deng et al. 2024; Liu et al. 2023c; Mu et al. 2021; Song et al. 2024; Swaminathan et al. 2024; Tseng et al. 2022; Wang et al. 2024b; Wei et al. 2022; Weng et al. 2024] and 3D Gaussian Splatting [Liu et al. 2025; Wu et al. 2025] to handle articulated objects. While effective at capturing object-specific geometry and appearance, these approaches are often slow due to their iterative optimization procedures. They also require densely sampled, posed images covering multiple articulation states, which are difficult to obtain in real-world settings, thereby limiting the scalability of these methods. Furthermore, they primarily rely on dense correspondence and lack cross-object priors, making it challenging to infer occluded or interior regions accurately.

2.3 Feedforward Generation of 3D Articulated Objects

Recent research has increasingly explored feedforward methods for the reconstruction and generation of 3D articulated objects, aiming to improve scalability and generalization. Some approaches focus on training 3D diffusion models specifically for articulated shapes [Gao et al. 2024; Lei et al. 2023; Luo et al. 2024], while others leverage fine-tuned vision-language models to predict joint parameters and part-level bounding boxes [Le et al. 2024]. Additionally, several methods adopt the paradigm of large-scale reconstruction model to synthesize novel views of articulated objects [Gao et al. 2025]. A diverse range of input modalities has been explored, including object categories and articulation graphs [Liu et al. 2024d], single-view images [Chen et al. 2024; Dai et al. 2024; Kawana and Harada 2023; Liu et al. 2024a], multi-view inputs [Gadi Patil et al. 2023; Heppert et al. 2023; Mandi et al. 2024; Zhang et al. 2021], textual descriptions [Su et al. 2024], and static 3D meshes [Qiu et al. 2025]. Although these feedforward approaches generally enable faster inference, their performance is constrained by the limited availability of large-scale datasets for 3D articulated objects [Xiang et al. 2020]. Consequently, they often rely on coarse geometric approximations—such as bounding boxes, predefined templates, or retrieval from small-scale databases—to represent object parts, limiting their ability to model fine-grained and realistic geometry. Moreover, most existing methods focus primarily on geometry reconstruction and overlook texture modeling, resulting in outputs with reduced visual realism. While a few recent approaches [Kawana and Harada 2023; Su et al. 2024; Zhang et al. 2021] employ implicit representations such as signed distance fields (SDFs) or occupancy fields to generate geometry, they still lack integrated texture synthesis and typically depend on brittle, multi-stage pipelines. In this work, we pursue a unified feedforward framework that jointly recovers detailed geometry, realistic textures, and accurate articulation structures for 3D articulated objects.

2.4 Feedforward 3D Static Object Generation

Open-world 3D generation of static objects has achieved notable progress. Recent approaches—such as CLAY [Zhang et al. 2024],

Trellis [Xiang et al. 2024], and Hunyuan3D [Zhao et al. 2025]—pioneer native 3D diffusion models trained directly in the 3D domain. Enabled by large-scale datasets [Deitke et al. 2023a,b] and advances in 3D representations, VAEs, and diffusion models, these methods produce accurate geometry, high-quality textures, and strong open-world generalization. However, extending this success to articulated objects by training native 3D diffusion models faces the critical challenge of limited data availability—current articulated datasets contain only a few thousand shapes [Xiang et al. 2020], which is orders of magnitude fewer than their static counterparts.

In parallel, many works investigate feedforward reconstruction models. Some directly regress 3D shapes from single-view images [Hong et al. 2023; Tochilkin et al. 2024; Zou et al. 2024], while others leverage sparse multi-view inputs [Li et al. 2023; Liu et al. 2024c, 2023d, 2024e; Tang et al. 2024; Wang et al. 2024a; Wei et al. 2024; Wu et al. 2024; Xu et al. 2024]. Though effective for static objects, these methods are rarely extended to articulated settings. Notably, a recent work, LVSM [Jin et al. 2024], demonstrates that a simple transformer model with minimal 3D inductive bias can achieve high-fidelity novel view synthesis for static objects. In this work, we build upon LVSM to support articulated objects—extending it to generate not only novel views under varying camera poses and joint states, but also explicit 3D textured meshes with disentangled parts and accurate articulation parameters.

3 Method

We introduce LARM, a unified framework for reconstructing 3D articulated objects from sparse-view RGB images captured at two distinct articulation states. The method enables articulation-aware novel view synthesis and supports explicit 3D reconstruction of both textured geometry and joint structure. We first describe the core model architecture, which encodes joint-aware input and target views to support both view synthesis and auxiliary predictions (Section 3.1). Next, we show how the outputs of LARM are used to estimate joint parameters and reconstruct explicit 3D textured meshes (Section 3.2). Finally, we outline a two-stage training strategy and data augmentation techniques that facilitate network convergence and improve generalization (Section 3.3).

3.1 LARM Model

Problem Formulation. LARM takes as input $N \times 2$ sparse-view images of the object captured at two different joint states, along with known camera extrinsics and intrinsics. The input can be represented as:

$$\{(\theta_i, I_{ij}, E_{ij}, K_{ij}), |, i = 1, 2; j = 1, \dots, N\}, \quad (1)$$

where θ_i denotes the articulation state, I_{ij} is the j -th input image at state θ_i , and E_{ij} and K_{ij} are the corresponding camera extrinsics and intrinsics. Given a target joint state θ_t and a novel camera pose defined by extrinsics E_t and intrinsics K_t , LARM synthesizes the corresponding target-view image I_t , along with optional outputs (e.g., depth maps or segmentation masks), which can be used for explicit 3D reconstruction.

While the two joint states, θ_1 and θ_2 , are used as input, they are not required to tie to specific physically interpretable quantities such as degrees or distances. For instance, one can simply assign the “rest” state as $\theta_1 = 0$ and the “maximum” joint state as $\theta_2 = 1$, regardless

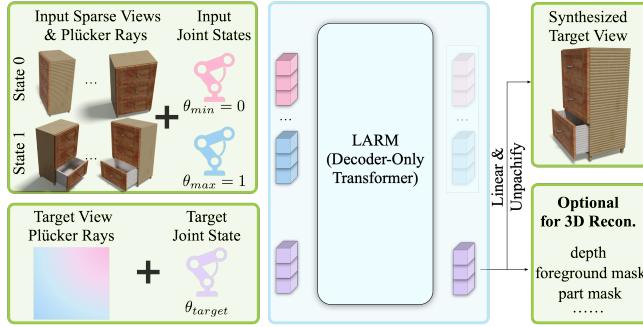


Fig. 2. LARM Architecture. LARM first patchifies the sparse, posed input images into tokens by concatenating the input RGB values, Plücker ray embeddings, and corresponding joint states. The target view to be synthesized is similarly represented by its Plücker ray embeddings and a target joint state, which are concatenated and tokenized. These input and target tokens are then fed into a decoder-only transformer model that predicts tokens used to regress the target view pixels. To enable explicit 3D reconstruction, LARM is also trained to produce additional outputs beyond RGB values, such as depth maps, foreground masks, and part masks.

of whether the actual articulation corresponds to 80 degrees or 0.3 meters. This relative scale avoids the need to annotate exact physical quantities for the input images. In this paper, we mainly focus on the cleaner and more useful setting where inputs only involve a single joint, such that θ_i is a scalar. However, in Section 4.7, we show that extending to multi-part inference, such as reconstructing a cabinet with multiple doors and drawers, is straight-forward and only requires minimal modifications to our pipeline.

Model Architecture As shown in Figure 2, LARM is built upon LVSM’s decoder-only architecture. For each input view I_{ij} , we first compute a pixel-wise Plücker ray embedding $P_{ij} \in \mathbb{R}^{H \times W \times 6}$ using the corresponding camera extrinsics and intrinsics, where H and W denote the height and width of the input images, respectively. Both the input image I_{ij} and its corresponding Plücker ray embedding P_{ij} are divided into patches of size $p \times p$, resulting in $I_{ijk} \in \mathbb{R}^{3p^2}$ and $P_{ijk} \in \mathbb{R}^{6p^2}$ for patch index $k = 1, \dots, HW/p^2$. For each patch, we concatenate the image patch, the Plücker ray embedding patch, and the corresponding joint state θ_i , yielding a vector $\text{concat}([I_{ijk}, P_{ijk}, \theta_i]) \in \mathbb{R}^{9p^2+1}$. This vector is then passed through a linear layer to produce an input patch token $x_{ijk} \in \mathbb{R}^d$, where d is the latent feature dimension. Target patch tokens are constructed similarly but only use the Plücker ray embedding patch from the target view P_{tk} and the target joint state θ_t . These are concatenated and mapped through another linear layer to obtain the target-view patch tokens $q_{tk} \in \mathbb{R}^d$.

We flatten and concatenate the token sequences from the input views, denoted as x_1, \dots, x_{l_x} , and from the target view, denoted as q_1, \dots, q_{l_q} , where $l_x = 2NHW/p^2$ and $l_q = HW/p^2$. The input tokens are fed into a decoder-only transformer model with multiple layers of self-attention. After applying the transformer layers, we obtain output tokens with the same sequence length as the input. Input view tokens x_1, \dots, x_{l_x} and target view tokens q_1, \dots, q_{l_q} are updated to x'_1, \dots, x'_{l_x} and q'_1, \dots, q'_{l_q} , respectively. Each output target token q'_i is passed through a linear layer to obtain a vector in

\mathbb{R}^{p^2c} , where c is the number of desired output channels. These vectors are reshaped into 2D patches, activated via a sigmoid function, and finally assembled to reconstruct the synthesized novel view \hat{I}_t .

Auxiliary Output and Loss Functions The original LVSM focuses solely on novel view synthesis of static objects and supervises training with the following loss:

$$\mathcal{L}_{RGB} = \text{MSE}(\hat{I}^t, I^t) + \lambda_{perceptual} \cdot \text{Perceptual}(\hat{I}^t, I^t), \quad (2)$$

where $\lambda_{perceptual}$ is a weight that balances the rendering loss and perceptual loss.

In our setup, we can use the same loss to supervise LARM for modeling variations across both camera poses and joint states. However, in applications such as building digital twins in robotics, novel view synthesis alone is insufficient—we also require explicit 3D reconstruction of object geometry and articulation. To address this, we extend the model to additionally predict depth maps \hat{D}_t and \hat{D}_{ij} , foreground masks \hat{MF}_t and \hat{MF}_{ij} , and movable part masks \hat{MP}_t and \hat{MP}_{ij} . Unlike LVSM, which discards input view tokens after transformer, LARM leverages both the updated input and target view tokens to generate these auxiliary outputs. Specifically, \hat{D}_t , \hat{MF}_t , and \hat{MP}_t are decoded from q'_1, \dots, q'_{l_q} , while \hat{D}_{ij} , \hat{MF}_{ij} , and \hat{MP}_{ij} are decoded from x'_1, \dots, x'_{l_x} , following the same decoding process used for synthesizing \hat{I}_t .

While only the auxiliary outputs corresponding to target views are used for explicit 3D reconstruction, the outputs from input views are leveraged during training to supervise the model and promote faster convergence. We supervise the full model using the following loss function:

$$\mathcal{L} = \mathcal{L}_{RGB} + \lambda_D \mathcal{L}_D + \lambda_{MF} \mathcal{L}_{MF} + \lambda_{MP} \mathcal{L}_{MP}, \quad (3)$$

where \mathcal{L}_D is the L1 loss applied to both predicted depth maps \hat{D}_{ij} and \hat{D}_t . \mathcal{L}_{MF} is the foreground mask loss for \hat{MF}_{ij} and \hat{MF}_t , computed as a binary cross-entropy (BCE) loss with foreground and background pixels averaged separately. \mathcal{L}_{MP} is the part mask loss for \hat{MP}_{ij} and \hat{MP}_t , and it consists of two terms. The first term is computed similarly to \mathcal{L}_{MF} , using a separately averaged BCE loss over the full image. The second term accounts for cases where the part occupies a small region of the image: we crop the bounding box around the part and apply the BCE loss to the cropped region. The coefficients λ_D , λ_{MF} , and λ_{MP} are the corresponding weights used to balance the loss terms.

3.2 Joint Estimation and Mesh Reconstruction

In this section, we describe how the outputs of the LARM model are used to estimate the underlying kinematic structure and reconstruct 3D meshes of articulated objects.

Joint Estimation. We consider two common joint types: prismatic and revolute, and take the joint type as input, which we regard as a low-barrier requirement.¹ For a revolute joint, we aim to estimate a pivot point $\mathbf{p} \in \mathbb{R}^3$ on the rotation axis, the axis direction $\mathbf{a} \in \mathbb{R}^3$, and a global joint scale factor s that maps joint state differences to actual angular displacements. For a prismatic joint, we similarly estimate the axis direction \mathbf{a} and the scale s .

To estimate these parameters, we query the LARM model with image pairs rendered under different joint states but with identical

¹For instance, using a single image of each object, ChatGPT correctly identified the joint types for 88.89% of our 144 test objects.

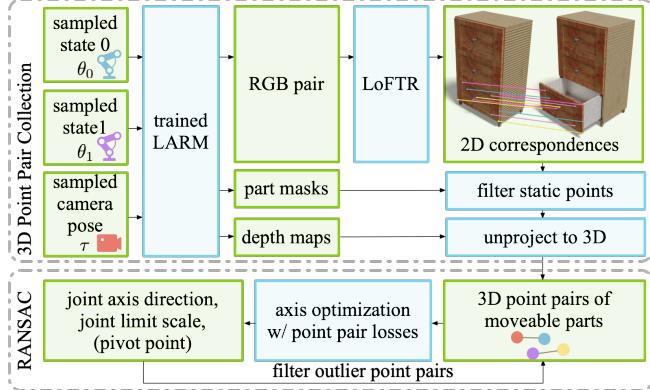


Fig. 3. Joint Estimation. To estimate explicit joint parameters using the LARM model, we first synthesize numerous image pairs with similar camera poses but different joint states. Next, we establish 2D pixel-wise correspondences for the movable part, which are then lifted to 3D. The joint parameters are optimized by minimizing the distances between the corresponding 3D point pairs under the estimated transformations. To enhance robustness, this optimization is integrated with RANSAC.

camera poses, as shown in Figure 3. Each pair provides a view of the object at two articulation states, θ_u and θ_v , while holding the viewpoint fixed. We then apply a 2D correspondence model (LoFTR [Sun et al. 2021]) to establish dense pixel-level matches between the two images. To ensure that the correspondences reflect true part motion, we retain only those matches located within the movable region, as defined by the predicted part masks \hat{M}_u and \hat{M}_v . We filter correspondence pairs based on LoFTR confidence and remove those that are too close in pixel space.

Given predicted depth maps and camera poses, the surviving 2D correspondences are unprojected to produce sets of 3D point pairs. Each pair (P_u^i, P_v^i) corresponds to a surface point observed at two joint states. The core idea is to model the transformation between these two points as a function of the joint parameters. For revolute joints, this transformation is a rotation around axis a passing through pivot p by an angle proportional to the joint state difference, $\Delta\theta = s(\theta_u - \theta_v)$. For prismatic joints, the transformation becomes a translation along a by a distance $\Delta d = s(\theta_u - \theta_v)$.

We then formulate an optimization problem to recover the joint parameters. The objective minimizes the 3D Euclidean distance between the transformed point P_u^i and its paired point P_v^i :

$$\mathcal{L}_{\text{joint}} = \sum_i \|T_{\theta_u \rightarrow \theta_v}(P_u^i; a, p, s) - P_v^i\|_2^2, \quad (4)$$

where $T_{\theta_u \rightarrow \theta_v}$ denotes the appropriate joint-induced transformation (either rotational or translational). We optimize this objective using gradient-based methods to jointly solve for a , p , and s . RANSAC [Fischler and Bolles 1981] is applied to eliminate outliers in joint estimation.

Mesh Reconstruction As shown in Figure 4, we sample a range of joint states and camera poses and use LARM to generate synthetic outputs for each configuration. The predicted depth maps are unprojected into 3D space using the known camera parameters, resulting in colored point clouds from multiple views. Each 3D point is associated with both a color (from the synthesized image) and a part label—either body or movable—based on the predicted part

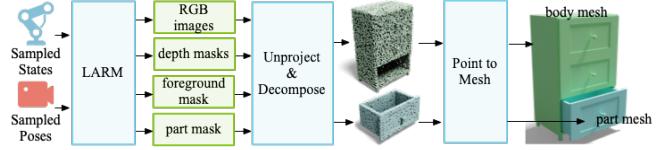


Fig. 4. Mesh Reconstruction. To extract an explicit mesh using the LARM model, we first synthesize multiple views from different camera poses. These views are lifted to 3D using the predicted depth and segmented based on the predicted masks, forming two separate colored point clouds. These point clouds are then fed into off-the-shelf point cloud-to-mesh tools [Peng et al. 2021] for explicit mesh reconstruction.

and foreground masks. This labeling enables the separation of raw point clouds into distinct semantic groups.

To reconstruct surface meshes, we first fuse the multi-view point clouds and then put the result into an off-the-shelf point cloud-to-mesh tool [Peng et al. 2021] for mesh reconstruction. Importantly, we perform this reconstruction process separately for the movable part and the main body, resulting in two disjoint meshes that can articulate independently based on the estimated joint parameters. Our method recovers a single set of canonical meshes and estimates joint parameters by leveraging multi-state observations. To obtain meshes under a specific articulation state, we simply transform the canonical meshes using the inferred joint parameters instead of reconstructing meshes again.

3.3 Data and Training Strategy

A major challenge in modeling articulated objects lies in the scarcity of high-quality datasets. To overcome this limitation and ensure both generalization and task-specific specialization, we adopt a two-stage training strategy: large-scale pretraining followed by task-specific finetuning.

Pretraining on Static Objects We begin by pretraining the LARM model on the Objaverse [Deitke et al. 2023b] dataset, which contains a diverse set of static 3D objects without articulation. During this stage, the model is trained purely for novel view synthesis using only RGB supervision \mathcal{L}_{RGB} . This stage equips the model with strong view synthesis capabilities and robust representations of object geometry and appearance. To further improve efficiency and performance, we adopt a coarse-to-fine training schedule: the model is first trained at a resolution of 256×256 and then refined at 512×512 resolution.

Finetuning on Articulated Objects After pretraining, we finetune the model on articulated objects from the PartNet-Mobility dataset, which includes detailed annotations of part articulation. In this stage, we activate the full prediction heads (for RGB, depth, masks, etc.) and supervise the model with the complete loss function introduced earlier. The linear output heads for RGB, depth, and masks are initialized by copying and repeating the pretrained RGB head weights, allowing the model to produce reasonable outputs from the very beginning of finetuning. This warm-start helps stabilize training and accelerates convergence.

Data Augmentation To enhance robustness and encourage generalization across unseen articulation patterns, we apply a series of augmentation strategies to the articulated dataset during finetuning. These include random scaling of objects along different axes to

Table 1. Novel View and State Synthesis Results. We report PSNR scores on our testing categories to reflect novel view synthesis quality, compared against Paris as the baseline method.

	Input	Storage.	Micro.	Refrig.	Safe	TrashCan	Table	Average
Paris	dense view	21.54	22.63	21.92	19.24	20.90	21.35	21.26
Ours	sparse view	33.06	30.42	31.25	30.06	29.89	29.89	30.76

Table 2. Cross-State Consistency. We generate videos of state-interpolation predictions and evaluate their temporal consistency by computing the Temporal LPIPS metric. This metric measures the LPIPS between adjacent frames and is computed over 1172 videos (25 frames each).

	Paris	Ours	GT
Temporal LPIPS	↓0.013	0.006	0.006

simulate variations in shape proportions, as well as random texture augmentation to diversify appearance.

4 Experiment

4.1 Implementation and Evaluation Details

In our experiments, we set $N = 3$, meaning that we use a total of 6 input images captured from two different joint states. The LARM model consists of 12 layers of self-attention with a bidirectional all-ones attention mask and a hidden dimension of 768. We first pretrain LARM on the novel view synthesis task using 3D static objects. Specifically, we train on a subset of 760k meshes from the Objaverse dataset [Deitke et al. 2023b]. For each mesh, we randomly sample 32 images to serve as the input and target view sources. The model is pretrained using 4 L40S GPUs for 3 days. After pretraining, we finetune the LARM model on articulated objects from the PartNet-Mobility dataset [Xiang et al. 2020]. We focus on six categories: StorageFurniture, Microwave, Refrigerator, Safe, TrashCan, and Table, with a total of 572 meshes (1,045 joints). We randomly hold out 116 shapes (with 293 joints in total) as test samples. To augment the training data, we apply random scaling and texture augmentation, where texture maps are randomly sampled from a database to replace the originals. For each joint movement, we generate eight variations (an object with two joints is augmented 16 times). For each variant, we render 32 views from randomly sampled camera poses and joint states, resulting in over 300,000 rendered views, which serve as the input and target images for training. Note that the train-validation-test splits are determined by objects rather than by images, which means all views of a test object are unseen during training. We finetune the LARM model using the loss described in Equation 3, employing 32 H100 GPUs for 5 days. The loss weights $\lambda_{\text{perceptual}}$, λ_D , λ_{MF} , and λ_{MP} are set to 0.1, 2.0, 1.0, 1.0 respectively. With the inferred novel views, mesh reconstruction and joint estimation for each object take approximately 90 seconds.

4.2 Novel View and State Synthesis

We compare our method with two recent state-of-the-art approaches, Paris [Liu et al. 2023c] and PartRM [Gao et al. 2025], on the novel view and state synthesis task using our test set from the PartNet-Mobility dataset [Xiang et al. 2020]. Paris performs NeRF-based optimization for each shape and requires dense-view input for two joint states. PartRM takes a single view and a 2D drag prompt as

input to synthesize a novel view after applying the drag operation. While both PartRM and our LARM model allow precise control over the joint state of the target view, we argue that the 2D drag prompt used by PartRM is less intuitive and less accurate for specifying the target joint state. For example, based on their released code, generating the drag prompt for evaluation on the PartNet-Mobility dataset requires access to the ground-truth 3D model. Although this setup is less practical, we still follow their protocol to generate the comparison results. Please refer to the supplementary materials for more details.

The qualitative results are shown in Figure 7. As observed, PartRM fails in many cases, generating twisted shapes. Paris produces better results but still often generates blurry outputs and floaters, particularly around the movable parts, despite utilizing dense-view inputs. We speculate that its per-shape NeRF optimization lacks cross-shape priors, making it less effective at handling challenging articulated components. In contrast, LARM demonstrates a stronger ability to synthesize novel views with sharp details and accurate articulated motion. The quantitative comparison in Table 1 further supports our observation, showing that LARM outperforms Paris by a large margin.

To evaluate the cross-state consistency of the synthesized views, we ask the model to generate images of continuous state interpolation and then assess their temporal consistency using temporal LPIPS, which measures the LPIPS between adjacent frames. As shown in Table 2, LARM maintains high temporal consistency, closely matching the ground truth and substantially outperforming the baseline Paris. We also provide qualitative examples in the supplementary materials.

4.3 3D Articulated Object Reconstruction

Beyond novel view and state synthesis, we compare our method with recent state-of-the-art approaches for 3D articulated object reconstruction: Articulate-Anything [Le et al. 2024], Singapo [Liu et al. 2024a], URDFFormer [Chen et al. 2024], and Paris [Liu et al. 2023c]. All methods except Paris operate on a single-view input. Articulate-Anything and Singapo construct articulated objects by retrieving parts and corresponding textures from the PartNet-Mobility dataset [Xiang et al. 2020], whereas URDFFormer predicts bounding boxes for individual parts, assembles template meshes, and projects the input image onto the mesh surface for texturing. During evaluation, we enumerate all joints of an object as target joints. For methods that predict multiple joints and links simultaneously, we compute the metric between their detected parts and the target part, selecting the one with the highest score. Please refer to the supplementary materials for more details.

We evaluate the reconstructed geometry (body and movable parts) based on both geometric and appearance accuracy. For each shape, we sample 5 joint states, generate corresponding meshes using the predicted joint parameters, and compute each metric for every state. The final results are reported as the average across all sampled states. Prior to evaluation, the predicted and ground-truth meshes are aligned using the method proposed in [Liu et al. 2024c]. For geometric evaluation, we compute Chamfer Distance (CD) [Fan et al. 2017] and F-Score [Wang et al. 2018], where CD is calculated using 100k sampled surface points, and F-Score is measured with

Table 3. Comparison of Reconstructed Mesh Geometry. We report Chamfer Distance and F-Score between the reconstructed meshes and ground-truth meshes from the PartNet-Mobility dataset. Each metric is first averaged across multiple joint states per object and then across different object instances. All baseline outputs are carefully aligned to match the format of PartNet-Mobility meshes, ensuring consistent and fair comparison. More detailed processing and setup of baseline methods can be found in the Appendix.

	Input	StorageFurniture		Microwave		Refrigerator		Safe		TrashCan		Table		Average	
		CD	F-Score	CD	F-Score	CD	F-Score	CD	F-Score	CD	F-Score	CD	F-Score	CD	F-Score
URDFFormer	single view	0.104	0.600	0.122	0.415	0.141	0.445	0.182	0.302	0.115	0.541	0.139	0.415	0.134	0.453
ArticulateAnything	single view	0.092	0.647	0.089	0.687	0.111	0.598	0.165	0.348	0.101	0.638	0.092	0.637	0.108	0.592
Singapo	single view	0.077	0.756	0.071	0.807	0.086	0.669	0.159	0.298	0.124	0.543	0.073	0.722	0.098	0.633
Paris	dense view	0.035	0.942	0.032	0.925	0.046	0.913	0.046	0.847	0.048	0.893	0.072	0.955	0.046	0.913
Ours	sparse view	0.023	0.960	0.028	0.931	0.034	0.908	0.046	0.835	0.026	0.961	0.020	0.977	0.030	0.929

Table 4. Comparison of Reconstructed Mesh Appearance. For each pair of reconstructed and ground-truth meshes, we render multi-view images and compute their PSNR and CLIP similarity. How we acquire textured mesh for comparison for each baseline method is detailed in the Appendix. All values are first averaged across multiple joint states and then across different shapes.

	StorageFurniture		Microwave		Refrigerator		Safe		TrashCan		Table		Average	
	PSNR	CLIP Sim	PSNR	CLIP Sim	PSNR	CLIP Sim	PSNR	CLIP Sim	PSNR	CLIP Sim	PSNR	CLIP Sim	PSNR	CLIP Sim
URDFFormer	15.1	0.886	12.3	0.893	14.2	0.893	12.2	0.886	14.3	0.861	12.8	0.846	13.5	0.878
ArticulateAnything	16.4	0.924	18.2	0.918	15.6	0.889	11.7	0.838	14.2	0.885	14.2	0.887	15.1	0.890
Singapo	11.7	0.843	9.9	0.886	13.3	0.887	11.1	0.871	11.5	0.826	11.2	0.873	11.4	0.864
Ours	22.5	0.919	22.5	0.906	23.3	0.912	20.1	0.902	22.0	0.906	20.2	0.903	21.8	0.908



Fig. 5. Real-world Demo. We use an iPhone to capture sparse-view images of everyday articulated objects. The results demonstrate that LARM can effectively handle such inputs and predict accurate novel views across diverse camera poses and joint states.

Table 5. Multi-Part Reconstruction. For each object, we sample five articulated instances by simultaneously sampling all joint states. When computing the metric, we first average over the five sampled articulated instances of each object and then average across objects.

	Input	StorageFurniture		Microwave		Refrigerator		Safe		TrashCan		Table		Average	
		CD	F-Score	CD	F-Score	CD	F-Score	CD	F-Score	CD	F-Score	CD	F-Score	CD	F-Score
Singapo	single view	0.104	0.603	0.131	0.504	0.113	0.509	0.158	0.381	0.149	0.471	0.120	0.517	0.114	0.556
Ours	3K+3 views	0.031	0.940	0.018	0.975	0.066	0.729	0.040	0.947	0.025	0.970	0.051	0.841	0.038	0.905

Table 6. Joint Estimation Accuracy. We report the success rates of four metrics; see the text for their definitions and thresholds.

Model	Axis Angle Error ↑	Axis Origin Error ↑	$M_r \uparrow$	$M_d \uparrow$
Ours	81.2%	81.9%	87.0%	84.6%
Singapo	75.1%	80.2%	35.2%	50.9%
ArticulateAnything	63.8%	64.5%	32.8%	42.3%
Paris	20.1%	58.4%	42.0%	37.2%
URDFFormer	51.9%	32.4%	27.7%	34.5%

a threshold of 0.05. For appearance evaluation, we render 8 random views for each state of both the predicted and ground-truth meshes, and compute the Peak Signal-to-Noise Ratio (PSNR) and CLIP similarity [Radford et al. 2021].

The qualitative results are shown in Figure 6. We observe that URDFFormer often produces overly simplistic geometry by compositing template board meshes and projecting the input images as

Table 7. **Ablation Study.** Evaluated on the table category.

	Model Variation	LARM Prediction				Mesh Reconstruction		
		Depth MAE	NVS PSNR	NVS LPIPS	mask mIoU	CD	F-Score	Mesh PSNR
a	w/o static objects pretraining	0.034	24.397	0.096	0.9309	—	—	—
b	w/o data augmentation	0.020	28.214	0.047	0.974	—	—	—
c	reduced number of views (64->16)	—	—	—	—	0.033	0.930	19.935
d	reduced number of qpos (5->2)	—	—	—	—	0.026	0.958	19.815
e	full	0.014	29.888	0.039	0.9832	0.020	0.977	20.194

textures, resulting in significant mismatches with the input prompts and ground-truth meshes. Articulate-Anything and Singapo achieve better results by retrieving part meshes from the PartNetMobility dataset. However, due to the limited scale of the retrieval database, they frequently capture only coarse geometry while failing to reproduce fine-grained geometric and texture details. In contrast, LARM, as a feedforward reconstruction model, faithfully reconstructs high-quality textured meshes that closely adhere to the input images.

Table 3 presents the quantitative evaluation of mesh geometry, where LARM outperforms all baselines across both metrics by a large margin. Notably, even though Paris utilizes dense-view inputs, it still performs slightly worse than our method, demonstrating the superiority of our feedforward model. Table 4 shows the quantitative evaluation of mesh appearance. Since the released code of Paris did not export mesh textures, we excluded it from the appearance comparison and only compared against the other three baselines. As shown in the table, our method outperforms all baselines by a large margin, particularly in the PSNR metric. The poor performance of the baseline methods is primarily due to their reliance on shape retrieval or template mesh compositing strategies, which lead to significant discrepancies between the reconstructed and ground-truth meshes.

4.4 Joint Parameter Estimation

We follow Articulate-Anything [Le et al. 2024] to evaluate joint parameter estimation accuracy using four metrics: (1) Axis angle error, the angle between the predicted and ground-truth axis; (2) Axis origin error, the shortest distance between the two axes; (3) Motion range distance M_r , which measures how far the predicted motion range deviates from the ground truth; (4) Motion direction difference M_d , which checks whether the predicted motion points the same way as the true motion. We report the success rate for these four metrics using thresholds of 0.25 radians for axis angle error, 0.15 for axis origin error, and 0.3 for both M_r and M_d . The results are shown in Table 6, where our method outperforms all baselines on all metrics, demonstrating more accurate joint estimation.

4.5 Ablation Study

We ablate our training and reconstruction strategies and report the results in Table 7.

Pretraining on 3D Static Objects. Due to the scarcity of 3D articulated objects, we first pretrain the model on a large-scale dataset of 3D static objects. This pretraining equips the model with fundamental knowledge of 3D object structures and the ability to synthesize novel views. When this pretraining is omitted and the model is trained from scratch on the limited set of articulated objects (row (a)), we observe a significant drop in performance, highlighting the importance of the pretraining stage.

Data Augmentation. To enrich the limited articulated objects, we apply augmentation of random scaling and the texture replacement. As shown in Row (b), this strategy effectively increases the data diversity and further improves the performance.

Number of Views for Mesh Reconstruction. LARM can synthesize novel views from arbitrary camera poses. For explicit mesh extraction, we use 64 views by default. In Row (c), we evaluate the effect of reducing this number to 16. While we observe a slight drop in performance, the results remain reasonable. Note that using fewer views facilitates faster reconstruction; however, when the number of views is further reduced to 4 or 8, holes occasionally appear in the reconstructed mesh.

Number of Joint States for Joint Estimation. LARM can synthesize outputs at novel joint states. To extract joint parameters (e.g., joint axis, pivot point), we query LARM at multiple joint states to construct correspondences and solve for the joint parameters. A greater number of joint states typically leads to more reliable correspondences and, consequently, more accurate joint estimation. In row (d), we reduce the default 5 sampled joint states to 2 and observe a drop in performance, highlighting the importance of LARM’s ability to synthesize outputs at novel joint states.

4.6 Real-world Evaluation

Beyond evaluating our method on standard synthetic datasets, we also assess its performance on real-world, sparsely captured images. As shown in Figure 5, LARM can effectively handle iPhone-captured images of everyday articulated objects. This demonstrates the broad applicability of our method and its potential to support various real-world applications, such as building digital twins.

4.7 Extension to Multi-Part Inference

Extending our method to multi-part objects requires only an inference-time adaptation. Here we assume that the object has multiple parts, but each part possesses only a single degree of freedom, such as a cabinet with multiple drawers and doors. For K movable parts, we use $3K + 3$ input images: three for the rest state (all parts closed) and three at the maximal articulation of each part. Each movable part is processed independently using its six images to infer the joint axis and reconstruct its part mesh, while the static body is recovered by fusing RGB-D point clouds and removing points belonging to movable parts via the inferred segmentation and depth.

We evaluate on the same test set under a multi-part setting, where the joint states of all joints are randomly sampled simultaneously and metrics are computed on the resulting articulated meshes. As shown in Table 5, our method remains effective in this multi-link scenario and outperforms the baseline Singapo. Note that we do not consider complex kinematic chains (e.g., a robot arm) in this work, which is consistent with most baselines that handle only single-DoF revolute or prismatic joints. This simplified setup covers a broad range of everyday articulated objects.

5 Conclusion and Limitation

We propose LARM, a novel feedforward model for reconstructing 3D articulated objects. Unlike recent approaches that rely on part retrieval or template meshes—which often result in significant mismatches compared to the input—LARM can reconstruct high-fidelity

novel views and 3D textured meshes that closely match the input prompts. Looking ahead, although we have leveraged pretraining on 3D static objects, the performance and generalizability of the model are still limited by the scale of the articulated object dataset. Exploring better strategies for data synthesis or more advanced training techniques may further improve the model's generalization and make it more widely applicable.

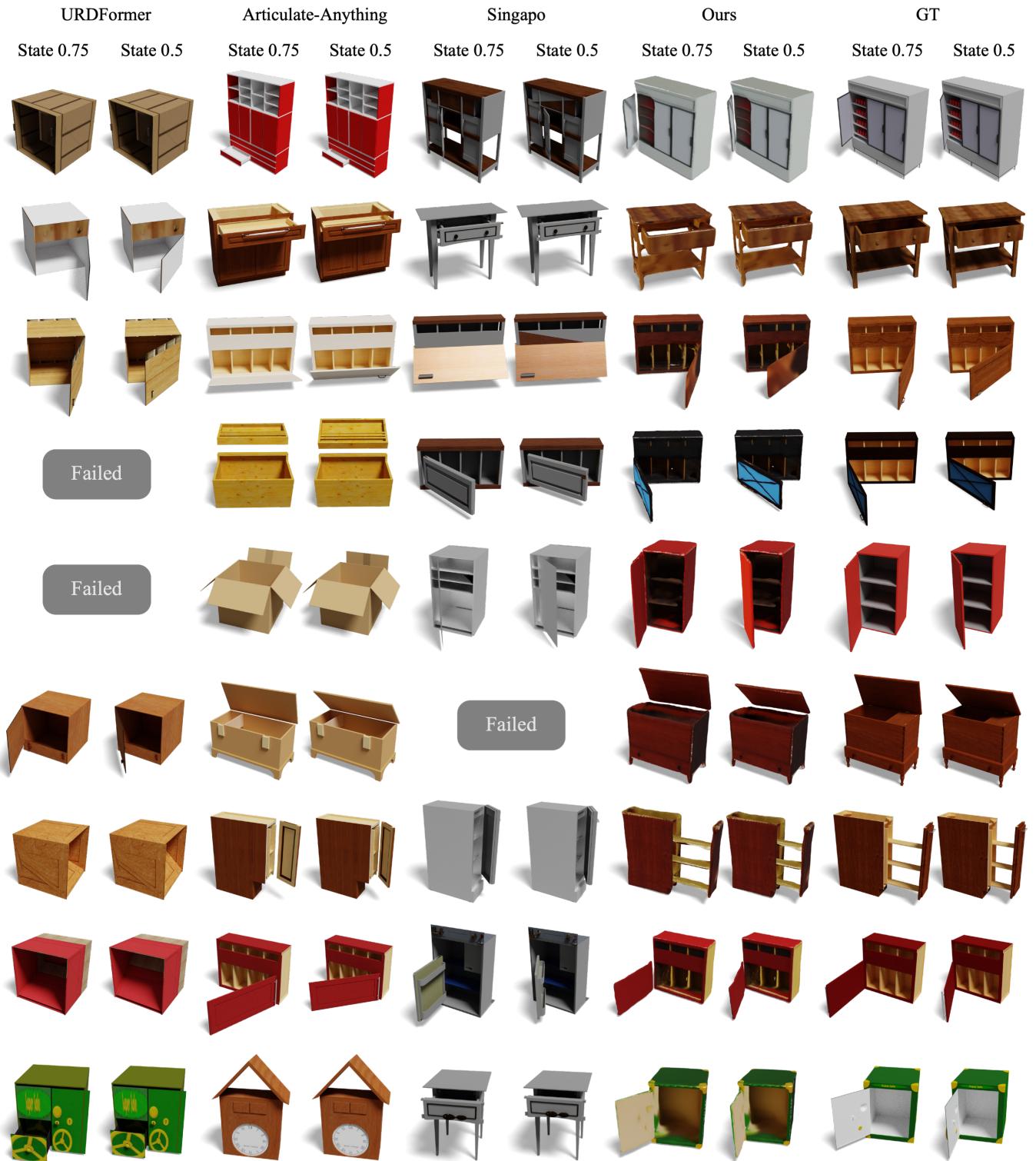


Fig. 6. Comparison of 3D Articulated Object Reconstruction. Note that methods such as URDFFormer [Chen et al. 2024] and Articulate-Anything [Le et al. 2024] rely on part retrieval for reconstruction, often resulting in significant mismatches in geometry and texture compared to the input prompt. In contrast, our LARM model faithfully reconstructs high-quality textured meshes that closely align with the input prompts.

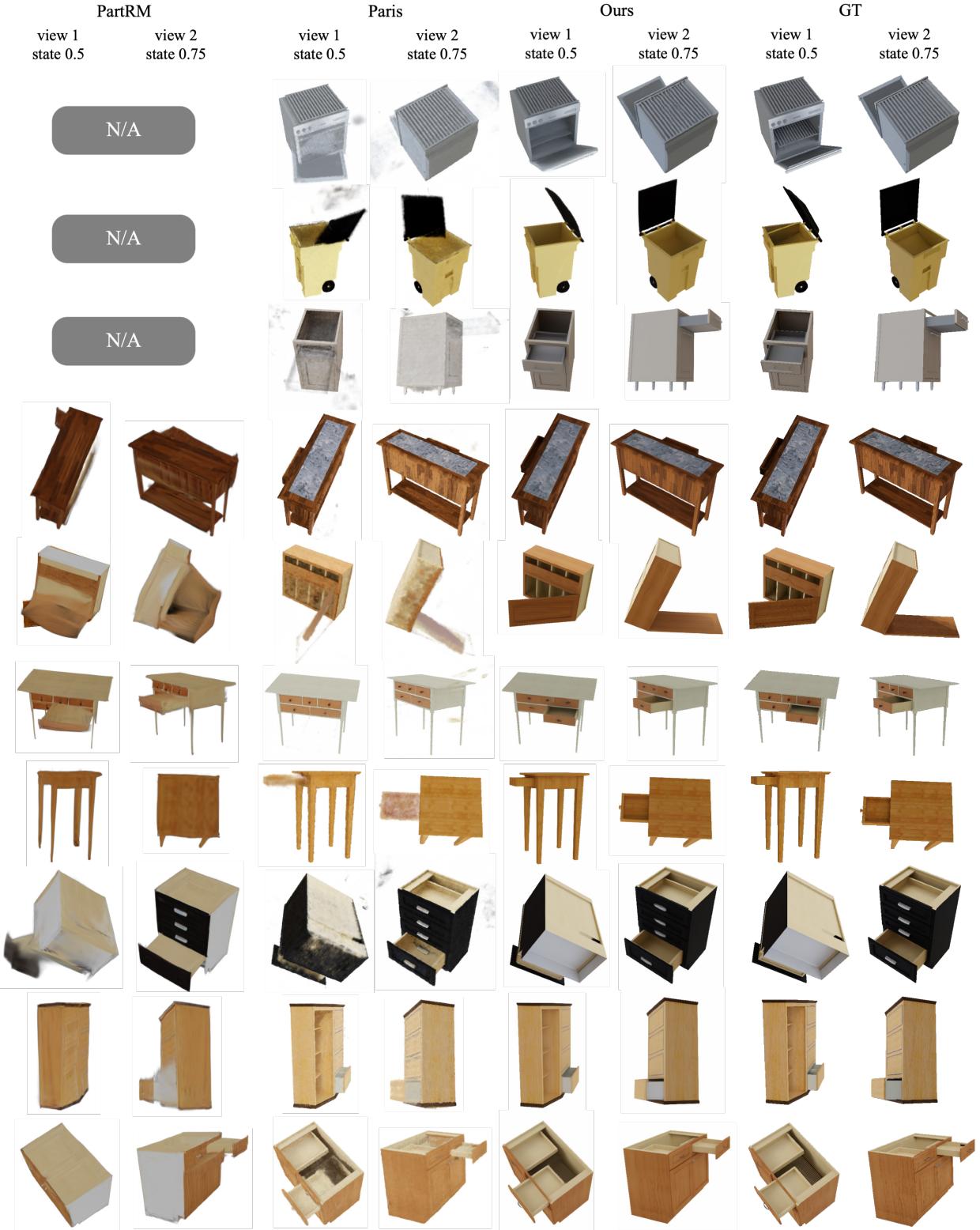


Fig. 7. Comparison of Novel View and State Synthesis between PartRM [Gao et al. 2025], Paris [Liu et al. 2023c], and our LARM. For each shape, we showcase synthesized views at two novel joint states. For objects not included in the PartRM preprocessed datasets, we leave blank. Please refer to the appendix for details on PartRM evaluation.

Appendix

A1 Baseline Setup and Evaluation Details

We evaluate each baseline under its as-published inference regime, using the released checkpoints when available or, for training-free methods, their provided prompts and inputs. Specifically, Articulate-Anything [Le et al. 2024] is a training-free, VLM-based method. For Paris [Liu et al. 2023c], there is no dataset-level training; instead, it performs per-scene optimization from multi-view images. For URDFFormer [Chen et al. 2024] and Singapo [Liu et al. 2024a], we use their released checkpoints, which are pretrained on their respective training splits. We did not retrain these models, as they may require additional supervision beyond our rendered images. All baseline methods and our method are evaluated on the same test set.

URDFFormer For evaluation of URDFFormer[Chen et al. 2024], we rendered and cropped a single front-view image from our Partnet-Mobility test set as input into the model, as a proxy to the real world image model input. We then used the provided fine-tuned Grounding DINO model for extracting predicted bounding boxes from the image. We did not manually adjust the output bounding boxes and instead directly used the predicted bounding boxes in the pipeline for articulation prediction. URDFFormer focuses on predicting articulation and joint information over producing realistic meshes and provides predefined mesh templates for visualization of output 3D meshes. In our evaluation, we used these predefined mesh templates, combined with their extracted texture from the input image, as the base meshes for rendering and computing various geometric and appearance metrics.

Articulate-Anything For evaluating Articulate-Anything[Le et al. 2024], we used its visual-input-based inference mode, providing single-view rendered images from our test set as input. We employed Google Gemini 1.5 Flash as the language model. Since the visual inference mode of Articulate-Anything is limited to retrieving entire objects from the PartNet-Mobility dataset, we excluded the target object itself from the candidate pool during inference to prevent exact matches and ensure a fair evaluation.

Singapo For evaluation of Singapo[Liu et al. 2024a], similarly, we provide a single-view rendering of our test set object as input. We use GPT-4o as the language model for acquiring input structural graphs from the input image. Singapo retrieves part meshes from partnet mobility and does not attempt to match the input image in texture. However, as a proxy, we instead enforce retrieval of texture along with any desired part mesh. The appearance and rendering metrics were computed with such retrieved textures.

Paris To construct the dense view input for Paris[Liu et al. 2023c], we rendered 100 training views for each input object joint, under the rest state and the maximum state respectively, ensuring same lighting condition as the target view renderings. Aside from novel view synthesis, we also used the extracted mesh for computation of geometry metrics.

For all of these above baseline methods, the resulting meshes are centered and normalized in the same way as the partnet mobility object meshes to ensure alignment and reasonable comparison of metrics.

Table 8. Category Generalization to an Unseen Class (Oven). LARM is trained without any oven instances and evaluated on the Oven category at test time. Each metric is first averaged across multiple joint states per object and then across object instances. Lower values indicate better performance for CD, whereas higher values are better for the other metrics.

Model	Oven				
	CD ↓	F-Score ↑	Render PSNR ↑	CLIP ↑	NVS PSNR ↑
URDFFormer	0.1335	0.4040	13.1878	0.8709	/
ArtAnything	0.0632	0.8146	14.4753	0.8818	/
Singapo	0.0714	0.7522	11.5511	0.8693	/
Paris	0.0343	0.9178	/	/	19.9293
Ours	0.0283	0.9272	20.5066	0.9671	28.0743

PartRM Because data preprocessing for PartRM[Gao et al. 2025] is particularly involved, we followed the exact evaluation pipeline from their released code— including drag prompts computed from the ground-truth meshes. To give PartRM the most favorable conditions, we bypassed their original Zero123++ stage (which infers multi-view images from a single view) and instead fed their reconstruction model the corresponding ground-truth images directly. Likewise, to eliminate any domain gap from differing camera intrinsics, we used the authors’ own preprocessed PartNet-Mobility as input and applied their pre-trained weights without alteration. Please note that these comparisons are strictly qualitative: our rendering parameters differ from theirs, and their preprocessed dataset omits several object categories. In those cases, we have left the corresponding slots blank in our visualizations.

A2 Qualitative Examples of Cross-State Consistency

In Figure 8, we showcase qualitative novel view synthesis results of articulation state interpolation, where our method exhibits strong cross-state consistency.

A3 Category Generalization

We acknowledge that our method’s generalizability is still limited by the scale of training data—a common challenge for feedforward approaches. However, unlike many baselines that rely on category-specific part structure templates or part retrieval databases, our method is both template-free and database-free, offering greater potential for generalization.

To evaluate the generalizability of LARM to unseen categories, we exclude *Oven* from training and assess performance on it at test time. As shown in Table 8, our model outperforms baselines trained with oven data across geometry (CD, F-Score), appearance (Render PSNR, CLIP), and novel-view synthesis (NVS PSNR), demonstrating stronger category-level generalization. While LARM generalizes well to categories similar to the training set, categories that differ substantially from the training distribution (e.g., scissors) remain challenging.

A4 Discussion of Failure Cases

We observe failures when (i) test shapes or kinematics deviate significantly from the training data (e.g., scissors), leading to inaccurate and blurry novel-view synthesis results, particularly in the movable part regions; (ii) sparse or self-occluded inputs undersample movable parts, causing unstable correspondences; and (iii) textureless or reflective regions impair depth prediction and 2D point matching. Figure 9 shows sample failure cases as discussed above.

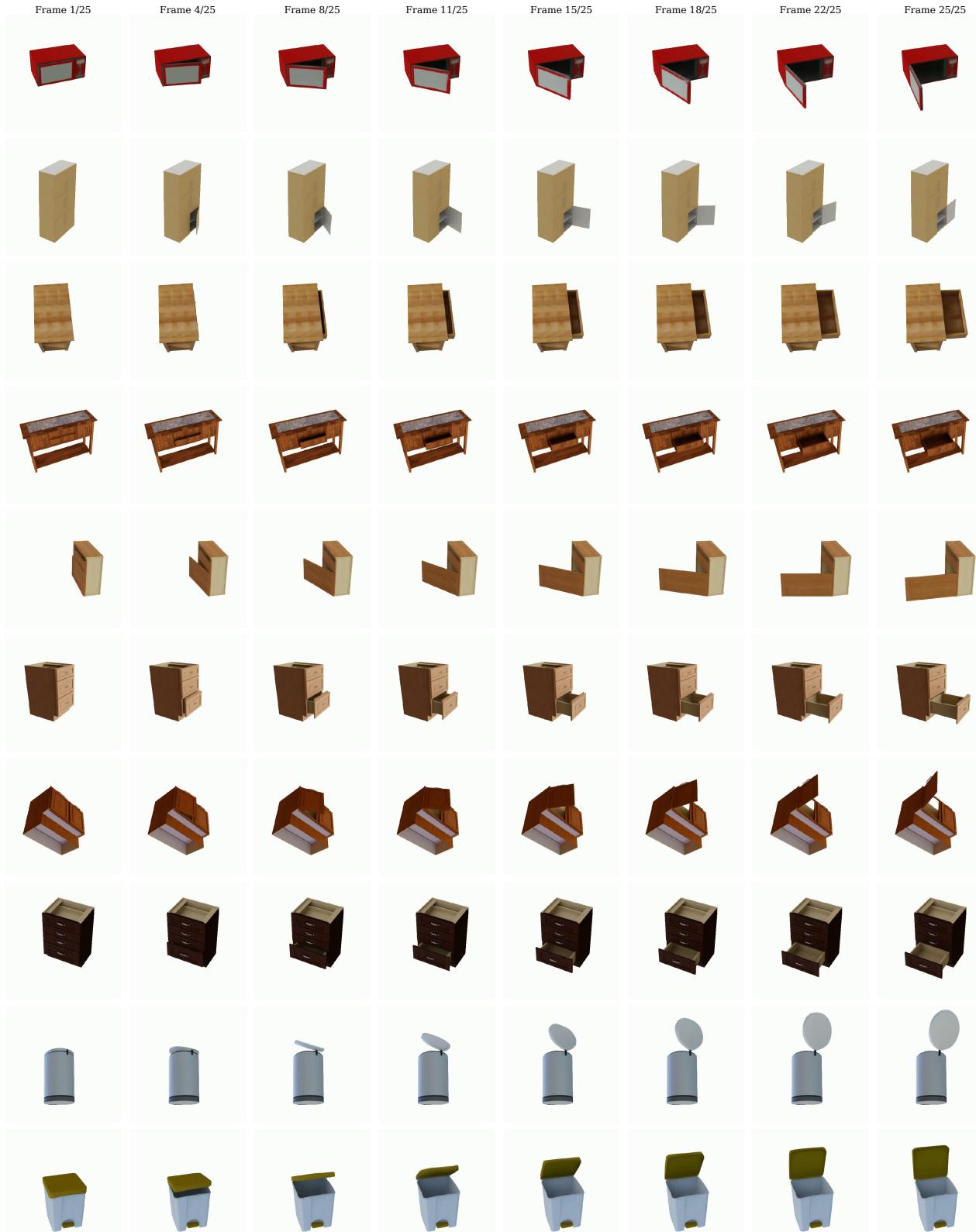


Fig. 8. **Qualitative Example of Cross-State Consistency.** For each example, we visualize NVS results across eight evenly sampled articulated states.



Fig. 9. **Failure Cases.** Here we show failure cases of LARM. Failure cases occur under three main categories: (i) test shapes are significantly different from the training data (e.g., scissors, rows 1 - 3); (ii) sparse or self-occluded inputs (rows 4-5); and (iii) large textureless or reflective regions on target objects (rows 6-7).

References

- Ben AbbateMatteo, Stefanie Tellex, and George Konidaris. 2019. Learning to generalize kinematic models to novel objects. In *Proceedings of the 3rd Conference on Robot Learning*.
- Yuchen Che, Ryo Furukawa, and Asako Kanezaki. 2024. OP-Align: Object-level and Part-level Alignment for Self-supervised Category-level Articulated Object Pose Estimation. In *European Conference on Computer Vision*. Springer, 72–88.
- Zoey Chen, Aaron Walsman, Marius Memmel, Kaichun Mo, Alex Fang, Karthikeya Vemuri, Alan Wu, Dieter Fox, and Abhishek Gupta. 2024. Urdformer: A pipeline for constructing articulated simulation environments from real-world images. *arXiv preprint arXiv:2405.11656* (2024).
- Tianyuan Dai, Josiah Wong, Yunfan Jiang, Chen Wang, Cem Gokmen, Ruohan Zhang, Jiajun Wu, and Li Fei-Fei. 2024. Automated creation of digital cousins for robust policy learning. *arXiv preprint arXiv:2410.07408* (2024).
- Matt Deitke, Ruoshi Liu, Matthew Wallingford, Huong Ngo, Oscar Michel, Aditya Kusupati, Alan Fan, Christian Laforet, Vikram Voleti, Samir Yitzhak Gadre, et al. 2023a. Objaverse-xl: A universe of 10m+ 3d objects. *Advances in Neural Information Processing Systems* 36 (2023), 35799–35813.
- Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihis, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. 2023b. Objaverse: A universe of annotated 3d objects. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 13142–13153.
- Jianning Deng, Kartic Subr, and Hakan Bilen. 2024. Articulate your NeRF: Unsupervised articulated object modeling via conditional view synthesis. *arXiv preprint arXiv:2406.16623* (2024).
- Haoqiang Fan, Hao Su, and Leonidas J Guibas. 2017. A point set generation network for 3d object reconstruction from a single image. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 605–613.
- Martin A. Fischler and Robert C. Bolles. 1981. Random Sample Consensus: A Paradigm for Model Fitting with Applications to Image Analysis and Automated Cartography. *Commun. ACM* 24, 6 (1981), 381–395. doi:10.1145/358669.358692
- Lian Fu, Ryoochi Ishikawa, Yoshihiro Sato, and Takeshi Oishi. 2024. CAPT: Category-level Articulation Estimation from a Single Point Cloud Using Transformer. *arXiv preprint arXiv:2402.17360* (2024).
- Akshay Gadi Patil, Yiming Qian, Shan Yang, Brian Jackson, Eric Bennett, and Hao Zhang. 2023. RoSI: Recovering 3D Shape Interiors from Few Articulation Images. *arXiv e-prints* (2023), arXiv–2304.
- Daoyi Gao, Yawar Siddiqui, Lei Li, and Angela Dai. 2024. MeshArt: Generating Articulated Meshes with Structure-guided Transformers. *arXiv preprint arXiv:2412.11596* (2024).
- Mingju Gao, Yike Pan, Huan-ang Gao, Zongzheng Zhang, Wenyi Li, Hao Dong, Hao Tang, Li Yi, and Hao Zhao. 2025. PartRM: Modeling Part-Level Dynamics with Large Cross-State Reconstruction Model. *arXiv preprint arXiv:2503.19913* (2025).
- Nick Heppert, Muhammad Zubair Irshad, Sergey Zakharov, Katherine Liu, Rares Andrei Ambrus, Jeannette Bohg, Abhinav Valada, and Thomas Kollar. 2023. Carto: Category and joint agnostic reconstruction of articulated objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 21201–21210.
- Nick Heppert, Toki Migitatsu, Brent Yi, Claire Chen, and Jeannette Bohg. 2022. Category-independent articulated object tracking with factor graphs. In *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 3800–3807.
- Yicong Hong, Kai Zhang, Juxiang Gu, Sai Bi, Yang Zhou, Difan Liu, Feng Liu, Kalyan Sunkavalli, Trung Bui, and Hao Tan. 2023. Lrn: Large reconstruction model for single image to 3d. *arXiv preprint arXiv:2311.04400* (2023).
- Ruižhen Hu, Wenchao Li, Oliver Van Kaick, Ariel Shamir, Hao Zhang, and Hui Huang. 2017. Learning to predict part mobility from a single static snapshot. *ACM Transactions On Graphics (TOG)* 36, 6 (2017), 1–13.
- Siyuan Huang, Haonan Chang, Yuhan Liu, Yimeng Zhu, Hao Dong, Peng Gao, Abdeslam Boularias, and Hongsheng Li. 2024. A3vlm: Actionable articulation-aware vision language model. *arXiv preprint arXiv:2406.07549* (2024).
- Ajinkya Jain, Stephen Giguere, Rudolf Lioutikov, and Scott Niekum. 2022. Distributional depth-based estimation of object articulation models. In *Conference on Robot Learning*. PMLR, 1611–1621.
- Ajinkya Jain, Rudolf Lioutikov, Caleb Chuck, and Scott Niekum. 2021. Screwnet: Category-independent articulation model estimation from depth images using screw theory. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 13670–13677.
- Hanxiao Jiang, Yongsen Mao, Manolis Savva, and Angel X Chang. 2022b. OPD: Single-view 3D Openable Part Detection. *arXiv preprint arXiv:2203.16421* (2022).
- Zhenyu Jiang, Cheng-Chun Hsu, and Yuke Zhu. 2022a. Ditto: Building digital twins of articulated objects from interaction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 5616–5626.
- Haian Jin, Hanwen Jiang, Hao Tan, Kai Zhang, Sai Bi, Tianyuan Zhang, Fujun Luan, Noah Shavely, and Zexiang Xu. 2024. Lvsrn: A large view synthesis model with minimal 3d inductive bias. *arXiv preprint arXiv:2410.17242* (2024).
- Yuki Kawana and Tatsuya Harada. 2023. Detection based part-level articulated object reconstruction from single RGBD image. *Advances in Neural Information Processing Systems* 36 (2023), 18444–18473.
- Long Le, Jason Xie, William Liang, Hung-Ju Wang, Yue Yang, Yecheng Jason Ma, Kyle Vedder, Arjun Krishna, Dinesh Jayaraman, and Eric Eaton. 2024. Articulate-anything: Automatic modeling of articulated objects via a vision-language foundation model. *arXiv preprint arXiv:2410.13882* (2024).
- Jiahui Lei, Congyue Deng, Bokui Shen, Leonidas Guibas, and Kostas Daniilidis. 2023. Nap: Neural 3d articulation prior. *arXiv preprint arXiv:2305.16315* (2023).
- Jiahua Li, Hao Tan, Kai Zhang, Zexiang Xu, Fujun Luan, Yinghao Xu, Yicong Hong, Kalyan Sunkavalli, Greg Shakhnarovich, and Sai Bi. 2023. Instant3d: Fast text-to-3d with sparse-view generation and large reconstruction model. *arXiv preprint arXiv:2311.06214* (2023).
- Xiaolong Li, He Wang, Li Yi, Leonidas J Guibas, A Lynn Abbott, and Shuran Song. 2020. Category-level articulated object pose estimation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 3706–3715.
- Jiayi Liu, Denys Iliash, Angel X Chang, Manolis Savva, and Ali Mahdavi-Amiri. 2024a. SINGAPO: Single Image Controlled Generation of Articulated Parts in Objects. *arXiv preprint arXiv:2410.16499* (2024).
- Jiayi Liu, Ali Mahdavi-Amiri, and Manolis Savva. 2023c. Paris: Part-level reconstruction and motion analysis for articulated objects. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 352–363.
- Jiayi Liu, Manolis Savva, and Ali Mahdavi-Amiri. 2024b. Survey on Modeling of Human-made Articulated Objects. *arXiv e-prints*, arXiv–2403.
- Jiayi Liu, Hou In Ivan Tam, Ali Mahdavi-Amiri, and Manolis Savva. 2024d. CAGE: controllable articulation generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 17880–17889.
- Liu Liu, Jianming Du, Hao Wu, Xun Yang, Zhenguang Liu, Richang Hong, and Meng Wang. 2023a. Category-level articulated object 9d pose estimation via reinforcement learning. In *Proceedings of the 31st ACM International Conference on Multimedia*, 728–736.
- Liu Liu, Han Xue, Wenqiang Xu, Haoyuan Fu, and Cewu Lu. 2022. Toward real-world category-level articulation pose estimation. *IEEE Transactions on Image Processing* 31 (2022), 1072–1083.
- Minghua Liu, Ruoxi Shi, Linghao Chen, Zhuoyang Zhang, Chao Xu, Xinyue Wei, Hansheng Chen, Chong Zeng, Jiayuan Gu, and Hao Su. 2024c. One-2-3-45++: Fast single image to 3d objects with consistent multi-view generation and 3d diffusion. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 10072–10083.
- Minghua Liu, Chao Xu, Haian Jin, Linghao Chen, Mukund Varma T, Zexiang Xu, and Hao Su. 2023d. One-2-3-45: Any single image to 3d mesh in 45 seconds without per-shape optimization. *Advances in Neural Information Processing Systems* 36 (2023), 2226–2226.
- Minghua Liu, Chong Zeng, Xinyue Wei, Ruoxi Shi, Linghao Chen, Chao Xu, Mengqi Zhang, Zhaoning Wang, Xiaoshuai Zhang, Isabella Liu, et al. 2024e. Meshformer: High-quality mesh generation with 3d-guided reconstruction model. *arXiv preprint arXiv:2408.10198* (2024).
- Qihao Liu, Weichao Qiu, Weiyao Wang, Gregory D Hager, and Alan L Yuille. 2020. Nothing but geometric constraints: A model-free method for articulated object pose estimation. *arXiv preprint arXiv:2012.00088* (2020).
- Shaowei Liu, Saurabh Gupta, and Shenlong Wang. 2023b. Building rearticulable models for arbitrary 3d objects from 4d point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 21138–21147.
- Yu Liu, Baoxiong Jia, Ruijie Lu, Junfeng Ni, Song-Chun Zhu, and Siyuan Huang. 2025. Building Interactable Replicas of Complex Articulated Objects via Gaussian Splatting. In *The Thirteenth International Conference on Learning Representations*.
- Rundong Luo, Haoran Geng, Congyue Deng, Puahao Li, Zan Wang, Baoxiong Jia, Leonidas Guibas, and Siyuan Huang. 2024. Physpart: Physically plausible part completion for interactable objects. *arXiv preprint arXiv:2408.13724* (2024).
- Zhao Mandi, Yijia Weng, Dominik Bauer, and Shuran Song. 2024. Real2code: Reconstruct articulated objects via code generation. *arXiv preprint arXiv:2406.08474* (2024).
- Jiteng Mu, Weichao Qiu, Adam Kortylewski, Alan Yuille, Nuno Vasconcelos, and Xiaolong Wang. 2021. A-sdf: Learning disentangled signed distance functions for articulated shape representation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 13001–13011.
- Songyou Peng, Chiyu Jiang, Yiyi Liao, Michael Niemeyer, Marc Pollefeys, and Andreas Geiger. 2021. Shape as points: A differentiable poisson solver. *Advances in Neural Information Processing Systems* 34 (2021), 13032–13044.
- Shengyi Qian, Linyi Jin, Chris Rockwell, Siyi Chen, and David F Fouhey. 2022. Understanding 3d object articulation in internet videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 1599–1609.
- Xiaowen Qiu, Jincheng Yang, Yian Wang, Zhehuan Chen, Yufei Wang, Tsun-Hsuan Wang, Zhou Xian, and Chuang Gan. 2025. Articulate Anything: Open-vocabulary 3D Articulated Object Generation. (2025).
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021.

- Learning transferable visual models from natural language supervision. In *International conference on machine learning*. PMLR, 8748–8763.
- Chaoyu Song, Jiacheng Wei, Chuan Sheng Foo, Guosheng Lin, and Faya Liu. 2024. Reacto: Reconstructing articulated objects from a single video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5384–5395.
- Jiayi Su, Youhe Feng, Zheng Li, Jinhua Song, Yangfan He, Botao Ren, and Botian Xu. 2024. Artformer: Controllable generation of diverse 3d articulated objects. *arXiv preprint arXiv:2412.07237* (2024).
- Jiaming Sun, Zehong Shen, Yuang Wang, Hujun Bao, and Xiaowei Zhou. 2021. LoFTR: Detector-free local feature matching with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 8922–8931.
- Xiaohao Sun, Hanxiao Jiang, Manolis Savva, and Angel Xuan Chang. 2023. OPDMulti: Openable Part Detection for Multiple Objects. *arXiv preprint arXiv:2303.14087* (2023).
- Archana Swaminathan, Anubhav Gupta, Kamal Gupta, Shishira R Maiya, Vatsal Agarwal, and Abhinav Shrivastava. 2024. Leia: Latent view-invariant embeddings for implicit 3d articulation. In *European Conference on Computer Vision*. Springer, 210–227.
- Jiaxiang Tang, Zhaoxi Chen, Xiaokang Chen, Tengfei Wang, Gang Zeng, and Ziwei Liu. 2024. Lgm: Large multi-view gaussian model for high-resolution 3d content creation. In *European Conference on Computer Vision*. Springer, 1–18.
- Dmitry Tochilkin, David Pankratz, Zexiang Liu, Zixuan Huang, Adam Letts, Yangguang Li, Ding Liang, Christian Laforte, Varun Jampani, and Yan-Pei Cao. 2024. Triposr: Fast 3d object reconstruction from a single image. *arXiv preprint arXiv:2403.02151* (2024).
- Wei-Cheng Tseng, Hung-Ju Liao, Lin Yen-Chen, and Min Sun. 2022. Cla-nerf: Category-level articulated neural radiance field. In *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 8454–8460.
- Haowen Wang, Zhen Zhao, Zhao Jin, Zhengping Che, Liang Qiao, Yakun Huang, Zhipeng Fan, Xiuquan Qiao, and Jian Tang. 2024b. SM 3: Self-supervised Multi-Task Modeling with Multi-view 2D Images for Articulated Objects. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 12492–12498.
- Nanyang Wang, Yinda Zhang, Zhuwen Li, Yanwei Fu, Wei Liu, and Yu-Gang Jiang. 2018. Pixel2mesh: Generating 3d mesh models from single rgb images. In *Proceedings of the European conference on computer vision (ECCV)*. 52–67.
- Xiaogang Wang, Bin Zhou, Yahao Shi, Xiaowu Chen, Qipeng Zhao, and Kai Xu. 2019. Shape2motion: Joint analysis of motion parts and attributes from 3d shapes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8876–8884.
- Zhengyi Wang, Yikai Wang, Yifei Chen, Chendong Xiang, Shuo Chen, Dajiang Yu, Chongxuan Li, Hang Su, and Jun Zhu. 2024a. Crm: Single image to 3d textured mesh with convolutional reconstruction model. In *European Conference on Computer Vision*. Springer, 57–74.
- Fangyin Wei, Rohan Chabra, Lingni Ma, Christoph Lassner, Michael Zollhöfer, Szymon Rusinkiewicz, Chris Sweeney, Richard Newcombe, and Mira Slavecheva. 2022. Self-supervised neural articulated shape and appearance models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 15816–15826.
- Xinyue Wei, Kai Zhang, Sai Bi, Hao Tan, Fujun Luan, Valentin Deschaintre, Kalyan Sunkavalli, Hao Su, and Zexiang Xu. 2024. MeshLRM: Large Reconstruction Model for High-Quality Meshes. *arXiv preprint arXiv:2404.12385* (2024).
- Yijia Weng, He Wang, Qiang Zhou, Yuzhe Qin, Yueqi Duan, Qingnan Fan, Baoquan Chen, Hao Su, and Leonidas J Guibas. 2021. Captra: Category-level pose tracking for rigid and articulated objects from point clouds. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 13209–13218.
- Yijia Weng, Bowen Wen, Jonathan Tremblay, Valts Blukis, Dieter Fox, Leonidas Guibas, and Stan Birchfield. 2024. Neural implicit representation for building digital twins of unknown articulated objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3141–3150.
- Di Wu, Liu Liu, Zhou Linli, Anran Huang, Liangtu Song, Qiaojun Yu, Qi Wu, and Cewu Lu. 2025. REArtGS: Reconstructing and Generating Articulated Objects via 3D Gaussian Splatting with Geometric and Motion Constraints. *arXiv preprint arXiv:2503.06677* (2025).
- Kailu Wu, Fangfu Liu, Zhihan Cai, Runjie Yan, Hanyang Wang, Yating Hu, Yueqi Duan, and Kaisheng Ma. 2024. Unique3d: High-quality and efficient 3d mesh generation from a single image. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Fanbo Xiang, Yuzhe Qin, Kaichun Mo, Yikuan Xia, Hao Zhu, Fangchen Liu, Minghua Liu, Hanxiao Jiang, Yifu Yuan, He Wang, et al. 2020. Sapien: A simulated part-based interactive environment. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 11097–11107.
- Jianfeng Xiang, Zelong Lv, Sicheng Xu, Yu Deng, Ruicheng Wang, Bowen Zhang, Dong Chen, Xin Tong, and Jiaolong Yang. 2024. Structured 3d latents for scalable and versatile 3d generation. *arXiv preprint arXiv:2412.01506* (2024).
- Jiale Xu, Weihao Cheng, Yiming Gao, Xintao Wang, Shenghua Gao, and Ying Shan. 2024. Instantmesh: Efficient 3d mesh generation from a single image with sparse-view large reconstruction models. *arXiv preprint arXiv:2404.07191* (2024).
- Xianghao Xu, Yifan Ruan, Srinath Sridhar, and Daniel Ritchie. 2022. Unsupervised kinematic motion detection for part-segmented 3d shape collections. In *ACM SIGGRAPH 2022 Conference Proceedings*. 1–9.
- Zihao Yan, Ruizhen Hu, Xingguang Yan, Luanmin Chen, Oliver Van Kaick, Hao Zhang, and Hui Huang. 2020. RPM-Net: recurrent prediction of motion and parts from point cloud. *arXiv preprint arXiv:2006.14865* (2020).
- Zihao Yan, Fubao Su, Mingyang Wang, Ruizhen Hu, Hao Zhang, and Hui Huang. 2023. Interaction-Driven Active 3D Reconstruction with Object Interiors. *arXiv preprint arXiv:2310.14700* (2023).
- Hongliang Zeng, Ping Zhang, Chengjiong Wu, Jiahua Wang, Tingyu Ye, and Fang Li. 2024. MARS: multimodal active robotic sensing for articulated characterization. *arXiv preprint arXiv:2407.01191* (2024).
- Ge Zhang, Or Litany, Srinath Sridhar, and Leonidas Guibas. 2021. Strobenet: Category-level multiview reconstruction of articulated objects. *arXiv preprint arXiv:2105.08016* (2021).
- Longwen Zhang, Ziyu Wang, Qixuan Zhang, Qiwei Qiu, Anqi Pang, Haoran Jiang, Wei Yang, Lan Xu, and Jingyi Yu. 2024. CLAY: A Controllable Large-scale Generative Model for Creating High-quality 3D Assets. *ACM Transactions on Graphics (TOG)* 43, 4 (2024), 1–20.
- Zibo Zhao, Zeqiang Lai, Qingxiang Lin, Yunfei Zhao, Haolin Liu, Shuhui Yang, Yifei Feng, Mingxin Yang, Sheng Zhang, Xianghui Yang, et al. 2025. Hunyuan3d 2.0: Scaling diffusion models for high resolution textured 3d assets generation. *arXiv preprint arXiv:2501.12202* (2025).
- Zi-Xin Zou, Zhipeng Yu, Yuan-Chen Guo, Yangguang Li, Ding Liang, Yan-Pei Cao, and Song-Hai Zhang. 2024. Triplane meets gaussian splatting: Fast and generalizable single-view 3d reconstruction with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 10324–10335.