

# CAST: Component-Aligned 3D Scene Reconstruction from an RGB Image

KAIXIN YAO\*, ShanghaiTech University, China and Deemos Technology Co., Ltd., China  
 LONGWEN ZHANG\*, ShanghaiTech University, China and Deemos Technology Co., Ltd., China  
 XINHAO YAN, ShanghaiTech University, China and Deemos Technology Co., Ltd., China  
 YAN ZENG, ShanghaiTech University, China and Deemos Technology Co., Ltd., China  
 QIXUAN ZHANG†, ShanghaiTech University, China and Deemos Technology Co., Ltd., China  
 WEI YANG, Huazhong University of Science and Technology, China  
 LAN XU‡, ShanghaiTech University, China  
 JIAYUAN GU‡, ShanghaiTech University, China  
 JINGYI YU‡, ShanghaiTech University, China



Fig. 1. CAST brings diverse 3D scenes to life from a single image, where the relations between objects shaped by their physical roles and interactions come together to form a cohesive and immersive virtual environment.

\*Equal contributions.

†Project leader.

‡Corresponding author.

Authors' addresses: Kaixin Yao, ShanghaiTech University, Shanghai, China and Deemos Technology Co., Ltd., Shanghai, China, yaokx2023@shanghaitech.edu.cn; Longwen Zhang, ShanghaiTech University, Shanghai, China and Deemos Technology Co., Ltd., Shanghai, China, zhanglw2@shanghaitech.edu.cn; Xinhao Yan, ShanghaiTech University, Shanghai, China and Deemos Technology Co., Ltd., Shanghai, China, yanxh@shanghaitech.edu.cn; Yan Zeng, ShanghaiTech University, Shanghai, China and Deemos Technology Co., Ltd., Shanghai, China, zengyan2024@shanghaitech.edu.cn; Qixuan Zhang, ShanghaiTech University, Shanghai, China and Deemos Technology Co., Ltd., Shanghai, China, zhangqx1@shanghaitech.edu.cn; Wei Yang, Huazhong University of Science and Technology, China, Shanghai, weiyang@hust.edu.cn; Lan Xu, ShanghaiTech University, China, Shanghai, xulan1@shanghaitech.edu.cn; Jiayuan Gu, ShanghaiTech University, China, Shanghai, gujy1@shanghaitech.edu.cn; Jingyi Yu, ShanghaiTech University, China, Shanghai, yujingyi@shanghaitech.edu.cn.

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 ACM 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.

© 2025 Association for Computing Machinery.

XXXX-XXXX/2025/5-ART \$15.00

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

Recovering high-quality 3D scenes from a single RGB image is a challenging task in computer graphics. Current methods often struggle with domain-specific limitations or low-quality object generation. To address these, we propose CAST (Component-Aligned 3D Scene Reconstruction from a Single RGB Image), a novel method for 3D scene reconstruction. CAST starts by extracting object-level 2D segmentation and relative depth information from the input image, followed by using a GPT-based model to analyze inter-object spatial relations. This enables understanding of how objects relate to each other within the scene, ensuring more coherent reconstruction. CAST then employs an occlusion-aware large-scale 3D generation model to independently generate each object's full geometry, using Masked Auto Encoder (MAE) and point cloud conditioning to mitigate the effects of occlusions and partial object information, ensuring accurate alignment with the source image's geometry and texture. To align each object with the scene, the alignment generation model computes the necessary transformations, allowing the generated meshes to be accurately placed and integrated into the scene's point cloud. Finally, CAST applies a physics-aware correction mechanism, which leverages a fine-grained relation graph to generate a constraint graph. This graph guides the optimization of object poses, ensuring physical consistency and spatial coherence. By utilizing Signed Distance Fields (SDF), the model effectively addresses issues such as occlusions, object penetration, and floating objects, ensuring that the generated scene accurately reflects real-world physical interactions. Experimental results

demonstrate that CAST significantly improves the quality of single-image 3D scene reconstruction, offering enhanced realism and accuracy in scene understanding and reconstruction tasks. CAST has practical applications in virtual content creation, such as immersive game environments and film production, where real-world setups can be seamlessly integrated into virtual landscapes. Additionally, CAST can be leveraged in robotics, enabling efficient real-to-simulation workflows and providing realistic, scalable simulation environments for robotic systems.

CCS Concepts: • Computing methodologies → Artificial intelligence.

Additional Key Words and Phrases: Open-Vocabulary Scene Reconstruction, Generative Pose Alignment, Occlusion-aware 3D Generation, Physical Consistency

## ACM Reference Format:

Kaixin Yao, Longwen Zhang, Xinhao Yan, Yan Zeng, Qixuan Zhang, Wei Yang, Lan Xu, Jiayuan Gu, and Jingyi Yu. 2025. CAST: Component-Aligned 3D Scene Reconstruction from an RGB Image. 1, 1 (May 2025), 19 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

Humans exist within clear networks of relations—family, friends, coworkers—that guide our decisions and behaviors. These connections shape our world and give it structure. Similarly, objects in a space also function within their own networks [Latour 2005], but less noticed. They do not just exist in isolation; their placement, design, and material arise from physical constraints, functional roles, and human design intentions and influence how we move, interact, and perceive space. For example, a chair leans against a table for support, a cup rests on a saucer, and a lamp’s light interacts with surrounding surfaces, casting shadows that shape the overall scene. Recognizing these relations is critical for accurate scene parsing, modeling, and, more recently, 3D generation, ensuring virtual environments feel as realistic and coherent as the real world.

Significant progress has been made in generating single objects from text or image prompts. Neural rendering approaches [Poole et al. 2022; Wang et al. 2024a] optimize implicit representations, while native 3D generators [Xiang et al. 2024; Zhang et al. 2023, 2024a] directly create 3D shapes and textures via end-to-end learning. While these methods show promise for individual objects, applying them to generate entire scenes by assembling objects sequentially faces with notable shortcomings. A key challenge is accurate pose estimation. Existing methods often assume objects are view-aligned, which is rarely the case in real-world scenes. Objects may appear in diverse orientations, constrained by design, physics, or partial occlusion. Yet, most existing methods prioritize geometric fidelity over pose alignment, leaving this critical aspect underexplored.

An even more fundamental issue arises from the lack of inter-object spatial relations. Even with accurate poses, generated scenes often suffer from physically implausible artifacts: objects penetrate one another, float, or fail to make contact where necessary. These errors stem from the absence of spatial and physical constraints that naturally bind objects together, much as human relations structure our social world. While some recent methods [Liu et al. 2022; Zhang et al. 2024b] encode spatial relations implicitly using encoder-decoder architectures, they remain limited to specific domains such as indoor scenes. Other scene-level generators [Dogaru et al. 2024]

position objects in a global coordinate system but neglect their relative poses and dependencies, further compromising realism and usability for downstream applications like editing, animation, and simulation.

To this end, we propose CAST, a Component-Aligned 3D Scene reconstruction method for compositional reconstruction of a 3D scene from a single RGB image. CAST generates high-quality 3D meshes for individual objects, along with their similarity transformations (rotation, translation, scale), ensuring alignment with the reference image and enforcing physically sound interdependencies. CAST starts by processing an unstructured RGB image using 2D foundation models (e.g., Florence-2 [Xiao et al. 2024], GroundingDINO [Liu et al. 2025], SAM [Ravi et al. 2024], Grounded-SAM [Ren et al. 2024]) to recognize, localize, and segment objects in an open-vocabulary manner. Off-the-shelf monocular depth estimators [Wang et al. 2024b] provide partial 3D point clouds and initial estimates of inter-object spatial relations, including relative transformations and scales.

The first core component of CAST is our perceptive 3D instance generator with two modules: an occlusion-aware object generation module and a pose alignment generation module. The object generation module employs a latent diffusion-based generative model to produce high-fidelity object meshes conditioned on partial image segments and optional point clouds. This module incorporates an occlusion-aware 2D image encoder capable of inferring occluded regions, ensuring robust feature extraction for image conditions. To improve robustness to real-world point cloud conditioning, we simulate partial point clouds with occluded regions during training, enabling the model to handle occlusion effectively. The pose alignment module features an alignment generative model that produces a transformed partial point cloud, aligning with the complete geometry implicitly represented in the latent space. The similarity transformation is derived from the generated transformed point cloud and the partial point cloud estimated from the camera. Unlike direct pose regression methods [Kehl et al. 2017; Labb   et al. 2020], our method estimates transformations through generation, capturing the multi-modal nature of pose alignment.

The second core component of CAST addresses inter-object spatial relations. Despite accurate pixel alignment, physically implausible artifacts such as penetration or floating can occur without explicit modeling of physical constraints. CAST introduces a physics-aware correction process to ensure spatial and physical coherence. GPT-4v [Achiam et al. 2023] is utilized to identify commonsense physical relations grounded in the input image, which are then used to optimize object poses based on these constraints. This process ensures that reconstructed scenes exhibit realistic physical interdependencies, making them suitable for applications like simulation, editing, and rendering.

Remarkably, CAST excels at generating perceptually realistic 3D scenes from a wide range of images, whether they are sourced from indoor or outdoor settings, real-world captured, or AI-generated. Unlike previous approaches [Dai et al. 2024; Liu et al. 2022], CAST supports open-vocabulary reconstruction, even for challenging, in-the-wild images, thanks to our deliberate pipeline design. Quantitatively, CAST surpasses strong baselines in the indoor dataset,

3D-Front [Fu et al. 2021], regarding object- and scene-level geometry quality. It also outperforms on perceptual and physical realism across a diverse set of images, including in-the-wild scenarios, as verified by visual-language models and user studies.

Given only a single image, CAST can faithfully reconstruct the scene, with detailed geometry, vivid textures of the objects, and more importantly, the spatial and physical interdependencies between them. This capability democratizes virtual creation: a single snapshot of a room or outdoor space becomes a fully realized 3D environment, where objects are precisely posed, interact naturally, and account for occlusions. Game developers can integrate real-world setups into immersive landscapes, and filmmakers can effortlessly generate intricate virtual sets—unlocking creative potentials. Beyond entertainment, CAST paves the way for smarter robots. It can facilitate the real-to-simulation pipeline [Li et al. 2024a; Torne et al. 2024] by enabling robotics researchers to construct digital replicas from real-world demonstration datasets with more efficient and scalable simulation workflows.

## 2 RELATED WORK

Transforming real-world scenes into the digital realm enhances our ability to understand, recreate, and interact with the 3D world around us. This practice is widely embraced in industries such as animation, film, gaming, architecture, and manufacturing. It enables the creation of immersive movie experiences, the digital preservation of historical relics, and the development of interactive environments for gaming. For example, James Cameron employed groundbreaking 3D scanning technology in *Avatar* (2009) to bring the lush, realistic world of Pandora to life. Similarly, in the gaming industry, *The Witcher 3: Wild Hunt* incorporated lifelike terrain and architectures inspired by real-world locations in Poland, blending authentic cultural and natural elements with imaginative, open-world exploration.

Photogrammetry is a widely used method to capture the physical world in high detail and translate it into digital form [Barron et al. 2021, 2022; Chen et al. 2018; Goesele et al. 2007; Kerbl et al. 2023; Mildenhall et al. 2020; Müller et al. 2022], but it requires tens to hundreds of images from multiple viewpoints, making it time-consuming, resource-intensive, and hard to scale. In contrast, single-image-based approaches are more efficient and scalable, requiring only one image that can be easily obtained from online repositories, eliminating the need for expensive scanning devices or multi-view setups.

### 2.1 Single Image Scene Reconstruction

Scene-level reconstruction from a single image presents challenges due to object diversity, occlusions, and the need to preserve spatial relations. A starting point is monocular depth estimation, where depth is inferred from a single image, typically generating a depth point cloud [Bhat et al. 2023; Piccinelli et al. 2024; Wang et al. 2024b; Yang et al. 2024b; Yin et al. 2023]. While it provides valuable information, it struggles with occlusions and hidden portions of the scene. To address this, novel view synthesis methods use representations like radiance fields [Tian et al. 2023; Yu et al. 2021, 2022] and 3D Gaussian [Szymanowicz et al. 2024a,b], learning occlusion priors

from 3D datasets [Chang et al. 2015; Dai et al. 2017; Geiger et al. 2013; Sun et al. 2018]. Despite these advances, monocular reconstruction methods still often struggle to provide detailed and precise scene representations.

Some methods focus on directly regressing geometries along with their semantic labels in the scene [Chen et al. 2024a; Chu et al. 2023; Dahnert et al. 2021; Gkioxari et al. 2022]. These approaches typically rely on scene datasets with ground truth object annotations, such as Matterport3D [Fu et al. 2021] and 3DFront [Fu et al. 2021], which are often small in scale and limited to indoor room environments. However, the feed-forward nature of these methods leads to the generation of geometries that often lack sufficient detail and quality.

To better film real world to digital, other methods turn to retrieval-based approaches [Dai et al. 2024; Gao et al. 2024b; Gümeli et al. 2022; Kuo et al. 2021; Langer et al. 2022], which enhance scene quality by searching for and replacing objects in a scene with similar objects from a pre-existing dataset. These methods incorporate advanced tools such as GPT-4 [Achiam et al. 2023], SAM [Kirillov et al. 2023; Ren et al. 2024], and depth priors to decompose scenes. While these methods improve scene realism by integrating real-world objects, they are constrained by the richness and scope of the datasets they rely on. For scenes outside the dataset’s domain, retrieval-based methods either produce erroneous results or fail to find suitable replacements, significantly degrading the quality of the reconstructed scene.

### 2.2 Reconstruction as Generation

With the continuous advancements in the field, the ability to create high-quality 3D digital assets from various types of open-vocabulary images or text prompts has significantly improved. This advancement has precipitated a paradigm shift where the single-view reconstruction problem evolves into a generative 3D synthesis framework. This paradigm change allows for the generation of 3D assets without being confined to a fixed dataset, enabling more flexible and scalable scene reconstruction.

Much of the current research in 3D asset generation focuses on distilling 3D geometry from 2D images generative models [Poole et al. 2022; Tang et al. 2023; Wang et al. 2024a]. More recent developments expand this approach by incorporating multi-view images for supervision [Liu et al. 2024, 2023c,a; Long et al. 2024; Voleti et al. 2025; Wu et al. 2024a], often trained on rendered images from large-scale object datasets like Objaverse [Deitke et al. 2023], to enhance view consistency during generation. Some approaches directly regress the shape and appearance of individual objects based on input image [Hong et al. 2023; Tang et al. 2025]. While these methods achieve satisfactory visual results, they frequently fail to reproduce fine geometric details. To improve the quality of 3D geometry, a growing body of work has moved away from 2D supervision entirely, opting instead to train directly on 3D assets [Deitke et al. 2024, 2023]. These methods produce high-quality object-level geometries with advanced processing techniques [Xiang et al. 2024; Zhang et al. 2023, 2024a]. However, such approaches focus on isolated objects and fail to address scene-level challenges, such as modeling spatial hierarchies, inter-object relations, and environment lighting. Scene generation remains underdeveloped due to the high computational and

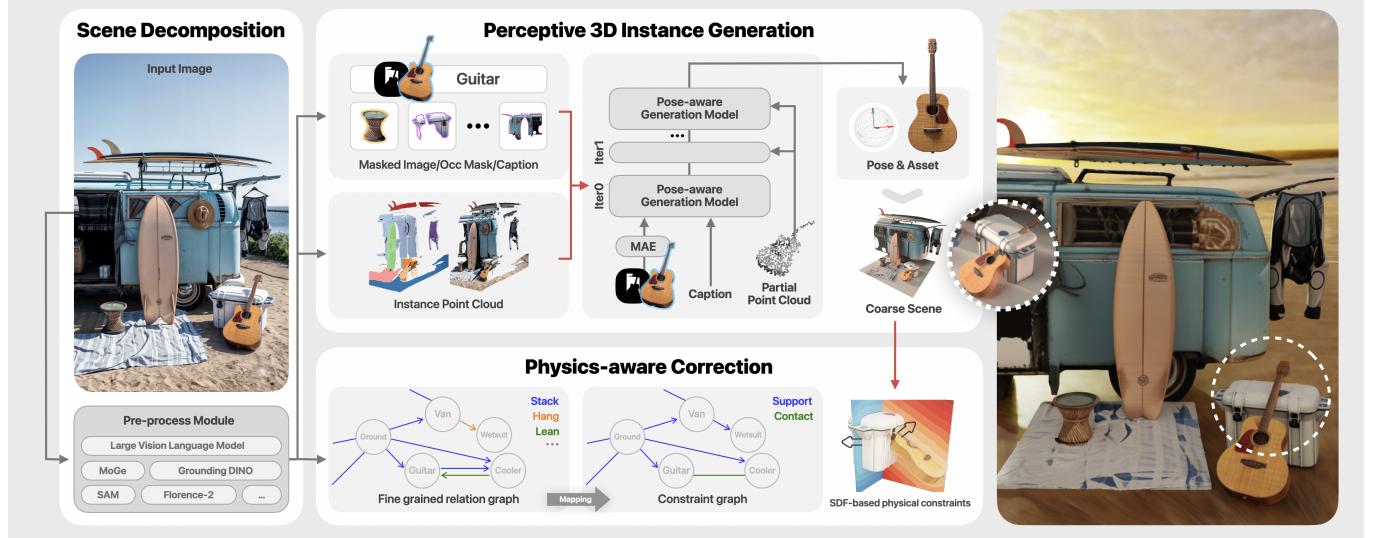


Fig. 2. Overview of the proposed pipeline. The input RGB image is processed through scene analysis to extract key information, followed by pose-aware generation to create initial 3D models. Physical constraint refinement ensures realistic interactions and spatial relations, yielding a high-quality, mesh-based 3D scene.

representational complexity of modeling object relations, lighting, and materials. Despite progress, current approaches still struggle to produce fully realized, editable 3D scenes. Existing paradigms either use video diffusion models [Blattmann et al. 2023; Ho et al. 2022a,b] to generate navigable 2D projections [Bruce et al. 2024; Yu et al. 2024], or rely on diffusion priors for volumetric scene approximations via 3D Gaussian splatting [Gao et al. 2024a; Liang et al. 2024; Wu et al. 2024b]. While these methods yield compelling visuals, they are incompatible with traditional production pipelines, lacking editable meshes, UV mappings, and decomposable PBR materials.

A more feasible paradigm decomposes scenes into modular components—objects, backgrounds, and environmental generating and reassembling them into an editable scene graph for greater flexibility and precision. For example, Gen3DSR [Dogaru et al. 2024] uses DreamGaussian [Tang et al. 2023] for open-vocabulary reconstruction. However, it struggles with occlusions, pose estimation, and editing individual objects, while relying on 2D models leads to poor geometric details and low-fidelity representations. Another recent work, Midi [Huang et al. 2024], learns spatial relations between objects in a scene but requires training on datasets with ground truth 3D meshes and annotations. This reliance on specific datasets limits its scalability and generalization to arbitrary scenes.

Our approach shares conceptual foundations with classical analysis-by-synthesis methods [Yuille and Kersten 2006], as both aim to infer 3D structure by generating explanations for observed imagery. However, while analysis-by-synthesis relies on iterative rendering and pixel-level optimization, our method leverages pre-trained generative models and learned priors to synthesize plausible 3D scenes directly, often bypassing explicit rendering and optimization loops, thereby improving scalability, efficiency, and adaptability to open-world scenarios.

Building on this foundation, we present a novel scene reconstruction pipeline that generates each object independently and aligns them into a cohesive scene. Unlike existing methods, our approach preserves accurate geometry, textures, and consistent spatial relations, resulting in more realistic, reliable, and editable reconstructions with improved quality and flexibility.

### 2.3 Physics-Aware 3D Modeling

Generating physically plausible 3D assets is crucial for ensuring realism and functionality in applications such as animation, gaming, and robotics. While recent 3D generative models excel at creating visually realistic objects, they often fall short in achieving physical plausibility. To address this limitation, physics-aware 3D generative models have been developed to integrate physical principles into the generation process. Some methods use soft-body simulation to animate 3D Gaussians [Xie et al. 2024; Zhong et al. 2025], or generate articulated objects with physics-based penalties [Liu et al. 2023b], while others ensure self-supporting structures through rigid-body simulation [Chen et al. 2024b; Mezghanni et al. 2022, 2021] or FEM [Guo et al. 2024; Xu et al. 2024]. These methods leverage offline or online physical simulations to check the physical validity of generated shapes and in turn guide generation. However, these approaches are typically confined to individual objects, overlooking the mutual influences between multiple objects within a scene.

Incorporating physical constraints into scene synthesis is much more challenging due to the inclusion of more complex relations, e.g., inter-object contact. Yang et al. [2024a] integrates constraints like object collisions, room layout, and object reachability into their scene-level generation pipeline. However, it is limited to indoor scene synthesis and relies on a closed-vocabulary database to perform shape retrieval. Ni et al. [2024] addresses the issue of physical implausibility in multi-view neural reconstruction. It leverages both

differentiable rendering and physical simulation to learn implicit representations. However, it requires multi-view images as input, focuses on individual objects, and primarily addresses only stability (simulating the dropping of objects). In contrast, our method operates in an open-vocabulary setting and requires only a single input image. Furthermore, it accounts for more complex inter-object relations, particularly support and contact, making it more versatile and applicable across diverse scenarios.

### 3 OVERVIEW

Scene-level reconstruction from a single image is a fundamental challenge in computer graphics, with broad applications in animation, virtual reality, and interactive gaming. Unlike object-level reconstruction, which focuses on isolated objects, scene-level reconstruction emphasizes the arrangement and relations of multiple entities under realistic (or stylized) physics. By capturing per-object structures, spatial relations, and contextual cues, this holistic approach enables more immersive experiences, compelling narratives, and efficient workflows—benefits that surpass those of single-object reconstructions. Although previous methodologies have explored feed-forward pipelines or retrieval-based approaches using fixed 3D templates [Dai et al. 2024; Liu et al. 2022], these methods often struggle to capture nuanced scene semantics and complex object relations. To address these limitations, we propose a generation-driven scene reconstruction approach with emphasized object relations to construct high-fidelity, contextually consistent 3D environments from a single, unannotated RGB image whether sourced from real-world photography or synthetic data (see Fig. 2).

A key insight of our method is the thorough object relation analysis of scene contextual information. First, we perform object segmentation to identify and localize constituent objects within the image. We then obtain preliminary geometric information, i.e., point clouds, and explore semantic and spatial relations among objects. This contextual backbone informs our subsequent object-wise generation pipeline, ensuring that each reconstructed object retains not only its geometric fidelity but also its correct placement in the broader scene. Finally, we synthesize a coherent 3D environment that respects physical plausibility—achieving structurally sound layouts and realistic interactions among scene elements.

Our research focuses on two primary objectives: to explore how generative models can effectively capture complex inter-object relations in order to produce realistic, scene-level reconstructions from a single image; and to identify strategies for integrating geometric cues and contextual information that maximize accuracy and plausibility in 3D reconstructions. Through this investigation, we demonstrate that generative methods provide a more flexible and robust alternative to traditional feed-forward and retrieval-based techniques. These methods allow for fine-grained control over both object-level details and global scene composition, thus streamlining content creation pipelines for animation, game development, and other fields requiring accurate, visually compelling 3D models. This work highlights the advantages of a generation-centric framework and lays the groundwork for future advancements in scene-level 3D reconstruction. It also underscores the growing importance of

context-driven approaches in bridging the gap between 2D imagery and immersive, interactive virtual environments.

*Preprocessing.* To facilitate comprehensive scene reconstruction from a single image, we first perform an extensive semantic extraction that provides a robust foundation for subsequent processing. Specifically, we employ Florence-2 [Xiao et al. 2024] to identify objects, generate their descriptions, and localize each object with bounding boxes. We then leverage GPT-4v [Achiam et al. 2023] to filter out spurious detections and isolate meaningful constituent objects, allowing for open-vocabulary object identification that is not constrained by predefined categories. Next, we use GroundedSAM-v2 [Ren et al. 2024] to produce a refined segmentation mask  $\{M_i\}$  for each labeled object  $\{o_i\}$ , thereby obtaining both precise object boundaries and corresponding occlusion masks, which play a crucial auxiliary role in the object generation stage. Apart from semantic cues, we also integrate geometric information by extracting a scene-level point cloud. Using MoGe [Wang et al. 2024b], we generate pixel-aligned point clouds  $\{q_i\}$  for each object  $\{o_i\}, i \in \{1, \dots, N\}$  and a global camera parameter in the scene coordinate system. This additional geometric data is subsequently matched to each object’s segmentation mask, providing a reliable structural reference for the final 3D scene reconstruction.

### 4 PERCEPTIVE 3D INSTANCE GENERATION

In the endeavor to reconstruct high-fidelity 3D scenes from single RGB images, a brute-force approach involves generating the entire scene mesh directly using techniques such as single-image depth estimation or diffusion priors. However, this method inherently struggles to manage occlusions, render invisible components, and accurately represent object relations due to the complex and intertwined nature of real-world scenes. Instead of generating an entire scene mesh directly, our approach focuses on individual object generation and then arranges the objects via precise relational alignment, as illustrated in 3. This strategy offers several advantages: 1. focusing on individual objects ensures higher geometric fidelity and allows for detailed modeling, resulting in more accurate and visually appealing scene components. 2. operating within a canonical space ensures that generated assets adhere to standardized orientations and scales, seamlessly integrating with artist-defined coordinate systems and promoting consistency across digital content creation tools. 3. the modular approach supports various applications such as editing, rendering, and simulation, enabling independent manipulation of objects for greater flexibility and efficiency. By decomposing scene reconstruction into object-wise generation and alignment, our method improves asset quality and manageability while enhancing the overall coherence and functionality of the 3D environment. This approach addresses challenges like geometric precision and efficient post-processing, advancing single-image 3D scene generation.

Object-wise generation presents significant challenges, primarily due to partial observations of objects within a scene caused by occlusions and limited sensor coverage. Additionally, existing generation methods often fail to coordinate multiple objects cohesively, resulting in inconsistent and unrealistic scenes. To overcome these limitations, we propose an Occlusion-Aware 3D Object Generation

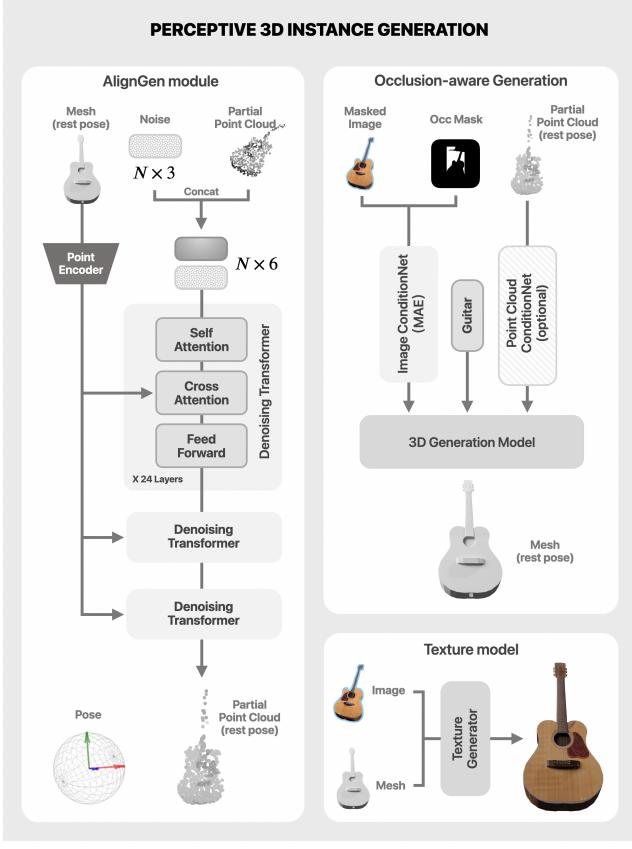


Fig. 3. Network design of our alignment generation model (Sec. 4.2), occlusion-aware object generation model (Sec. 4.1), and an illustrative figure of the texture generation model.

framework that integrates partial observations with comprehensive scene understanding. Specifically, given an image and its point cloud, our framework generates a high-quality 3D asset that not only resembles the input image but also aligns accurately with the partial point cloud represented in its corresponding canonical space. Furthermore, we compute a transformation matrix that maps the generated object from its canonical space back to the original scene space, ensuring spatial consistency within the scene.

A critical aspect of our object generation process is the utilization of a large generative model to generate holistic and high-fidelity object meshes from partial image and point cloud observations. To do so, we first follow state-of-art native 3D generative models [Xiang et al. 2024; Zhang et al. 2023, 2024a] to pre-train a large-scale 3D generative model conditioning on textual and image inputs.

We build upon existing generative frameworks featuring 3DShape2Vec representation [Zhang et al. 2023, 2024a], which prioritize geometry generation by utilizing a Geometry Variational Autoencoder (VAE). This VAE framework encodes uniformly sampled surface point clouds into unordered latent codes and decodes these latent representations into Signed Distance Fields (SDFs). Formally, the

VAE encoder  $\mathcal{E}$  and decoder  $\mathcal{D}$  are defined as:

$$Z = \mathcal{E}(X), \quad \mathcal{D}(Z, p) = SDF(p), \quad (1)$$

where  $X$  represents the sampled surface point cloud of the geometry,  $Z$  is the corresponding latent code, and  $SDF(p)$  denotes the operation of querying the SDF value at point  $p$  for subsequent mesh extraction via marching cubes. To effectively incorporate image information into the geometry generation process, we employ DINOv2 [Oquab et al. 2023] as our image encoder, following methodologies outlined in Xiang et al. [2024]; Zhang et al. [2023, 2024a]. The geometry latent diffusion model (LDM) is then formulated as:

$$\epsilon_{\text{obj}}(Z_t; t, c) \rightarrow Z, \quad (2)$$

where  $\epsilon$  represents the diffusion transformer model,  $Z_t$  is noisy geometry latent code at timestep  $t$ , and  $c$  denotes the encoded image features from DINOv2. We follow the pre-training process of prior works [Zhang et al. 2023, 2024a] and pre-train the base model on Objaverse [Deitke et al. 2023]. Upon training, our generation model  $\epsilon$  is capable of generating detailed 3D geometry solely based on image features.

#### 4.1 Occlusion-aware 3D Object Generation

Directly applying 3D generative-based models faces considerable challenges as real-world scenarios often present challenges such as partial occlusions in the input images, which severely degrade the quality and accuracy of the generated object geometries. To address this issue, we leverage the Masked Auto Encoder (MAE) capabilities of DINOv2. Specifically, during inference, we provide an occlusion mask  $M$  alongside the input image  $I$ , enabling the encoder to handle missing pixels by inferring latent features for the occluded regions. This is formalized as:

$$c_m = \mathcal{E}_{\text{DINOv2}}(I \odot M), \quad (3)$$

where  $M$  is a binary mask indicating which tokens should be masked and replaced with a [mask] token. During the pretraining phase, DINOv2 is trained with randomly set masks, allowing it to robustly infer missing parts based on the visible regions. Consequently, during inference, even if parts of the object image are occluded, the encoder can effectively reconstruct the necessary features, ensuring that the generative model maintains high-quality and accurate 3D reconstructions. This integration of image conditioning and occlusion handling is pivotal for our pipeline, as it ensures that the generated 3D objects are both visually consistent with the input images and geometrically faithful to the underlying structure.

*Canonical Point Cloud Conditioning.* Though our object generation model produces visually plausible meshes from input object images, it is challenging to generate pixel-aligned geometry due to the high-level nature of the encoded image condition  $c$  and the absence of pixel-wise supervision. We address this issue by additionally conditioning our object generation model on observed partial point clouds in canonical coordinates. This dual conditioning ensures that the generated geometries not only align visually with the input images but also accurately reflect their underlying scale, shape, and depth. During the conditioning training, we simulate real-world partial scans or estimated depth maps by rendering each 3D asset from multiple viewpoints, thereby obtaining corresponding RGB

images, camera parameters, and ground-truth depth maps. These RGB images are then processed using advanced depth estimation techniques, including MoGe [Wang et al. 2024b] and Metric3D [Yang et al. 2024b], to produce an estimated depth map and then projected as partial point clouds. To ensure scale consistency, we align estimated depth maps from MoGe and Metric3D with ground-truth depth maps by scaling and shifting them based on the median and median absolute deviation of valid depth values. The resulting point clouds are then normalized to a canonical  $[-1, 1]^3$  space to ensure consistent spatial representation for coarse object alignment.

To bolster the model’s robustness and its ability to generalize across diverse real-world scenarios, we employ a data augmentation strategy that interpolates between ground-truth partial point clouds  $\mathbf{p}_{\text{gt}}$  (projected from ground truth depth map to simulate accurate depth) and noisier, estimated partial point clouds  $\mathbf{p}_{\text{est}}$  (projected from estimated depth map and aligned to simulate estimated noisy depth from RGB). This interpolation is mathematically represented as:  $\mathbf{p}_{\text{disturb}} = \alpha \cdot \mathbf{p}_{\text{gt}} + (1 - \alpha) \cdot \mathbf{p}_{\text{est}}$ , where  $\alpha \in [0, 1]$  is a weighting factor which is sampled uniformly during training. Our object generator, named “*ObjectGen*”, with partial point cloud conditioning is formulated as

$$\epsilon(Z_t; t, \mathbf{c}, \mathbf{p}_{\text{disturb}}) \rightarrow Z, \quad (4)$$

where the conditioning adaptation scheme is based on attention mechanism similar to Zhang et al. [2023, 2024a]. Additionally, to mimic real-world occlusions and missing data, we randomly mask sets of basic primitives—such as circles and rectangles—in the depth maps from various camera views. This results in partial point clouds with occluded and incomplete regions, further enhancing the model’s ability to handle imperfect inputs. A critical design choice in our approach is to maintain the alignment of partial point clouds with the geometry in our training data set. Unlike methods that apply random scaling, translation, or rotation to augmented point clouds, our aligned partial point clouds ensure that the generative model can more effectively conform to the input point clouds’s inherent structure. This alignment restricts the model to adhere closely to the actual shapes and scales of objects, thereby facilitating more precise and coherent 3D reconstructions. By conditioning on these well-aligned partial point clouds, our model achieves superior alignment both in overall size and local geometric details, resulting in high-quality and reliable 3D geometry generation.

## 4.2 Generative Alignment

Each generated 3D object is within a normalized volume and assumes a canonical pose that may not be aligned with the image and scene space point cloud. This is because the image conditions use high-level features, such as DINOv2, to achieve better generalization. Ensuring that each object is correctly transformed and scaled to align with its presentation in the scene is crucial for scene composition. Though traditional alignment methods, such as Iterative Closest Point (ICP) [Arun et al. 1987; Best 1992], can be employed, they often fail to account for semantic context, leading to frequent misalignments and diminished accuracy (see Fig. 9). Instead, we introduce an alignment generative model conditioned on the scene-space partial point cloud  $\mathbf{q} \in \mathbb{R}^{N \times 3}$  and the canonical-space geometry latent code  $Z$ . Formally, we define our alignment

generator “*AlignGen*” as:

$$\epsilon_{\text{align}}(\mathbf{p}_t; t, \mathbf{q}, Z) \rightarrow \mathbf{p}, \quad (5)$$

where  $\epsilon_{\text{align}}$  is a point cloud diffusion transformer,  $\mathbf{p} \in \mathbb{R}^{N \times 3}$  is the transformed version of the scene-space partial point cloud to the canonical space, aligning with the generated object mesh.  $Z$  is the generated geometry latent of object corresponding to  $\mathbf{p}$  from the object generation model.  $\mathbf{p}_t$  is the noised version of  $\mathbf{p}$  at timestep  $t$ . In essence, the generation model maps the scene-space partial point cloud  $\mathbf{q}$  to  $\mathbf{p}$  in the canonical  $[-1, 1]^3$  space, aligning it with the generated object mesh. We can subsequently recover the similarity transformation (i.e., scaling, rotation, and translation) from  $\mathbf{q}$  and  $\mathbf{p}$  using the Umeyama algorithm [Umeyama 1991] as they are point-wise corresponded. This final step is numerically more stable than directly predicting transformation parameters.

In practice, we employ distinct conditioning strategies for the input point cloud  $\mathbf{q}$  and the geometric latent  $Z$ . For  $\mathbf{q}$ , we concatenate the input point cloud with the diffusion sample  $\mathbf{p}_t$  along the feature channel dimension, enabling the transformer architecture to learn explicit correspondences between the noisy canonical-frame partial cloud and the world-space partial cloud. For the geometric latent  $Z$ , we apply a cross-attention mechanism to inject it into the point diffusion transformer. This approach ensures that the model effectively incorporates spatial and geometric relations. Additionally, due to symmetry and replicated geometrical shapes, multiple valid  $\mathbf{p}$  may exist for a given  $\mathbf{q}$  and  $Z$ . Our diffusion model addresses this by sampling multiple noise realizations and aggregating the resulting transformations, to select the most confident and coherent representations.

## 4.3 Iterative Generation Procedure

Recall that in our design, the object point cloud is unusable for object generation initially, as it is represented in the scene space, while our object generation model requires canonical-space point cloud for conditioning. Solely depending on image cues for object generation often fails to produce pixel-aligned geometry, mainly because of the high-level semantic conditioning and biases inherent in 3D datasets. Fortunately, our design enables seamless integration of the object generation and alignment modules through a joint, iterative process. This integration ensures that each generated 3D object is not only visually consistent with the input image but also accurately positioned and scaled within the scene. The iterative workflow with step index  $k$  can be summarized in three key steps:

*Step 1: Object Generation.* For an object image with mask, the **Object Generation** module (Sec. 4.1) synthesizes the geometry latent code  $\mathbf{z}^{(k)}$  based on the image features  $\mathbf{c}$  derived from DINOv2 and the aligned point cloud  $\mathbf{p}^{(k)}$  in canonical coordinates. We set  $\mathbf{p}^{(0)}$  to scene space point cloud  $\mathbf{q}$  and set the point cloud conditioning scale factor  $\beta^{(k)}$  to progressively increase from 0 to 1 as iteration procedure goes on, allowing the partial point cloud to take influence over time. Formally, this process is represented as:

$$\mathbf{z}^{(k)} = \text{ObjectGen}(\mathbf{c}, \mathbf{p}^{(k)} \otimes \beta^{(k)}). \quad (6)$$

Hence our object generator solely relies on masked image conditioning in the first step. The latent code  $z^{(k)}$  is then decoded into a 3D geometry using the VAE decoder  $\mathcal{D}$ .

*Step 2: Alignment.* Subsequently, the **Generative Alignment** module (Sec. 4.2) takes the newly generated geometry latent code  $z^{(k)}$  and the partial point cloud  $q$  in scene coordinates to predict a transformed canonical-space partial point cloud  $p^{(k+1)}$ :

$$p^{(k+1)} = \text{AlignGen}(q, z^{(k)}). \quad (7)$$

This transformed point cloud  $p^{(k+1)}$  serves as an improved alignment reference for the next iteration. By leveraging the generative transformation model, the model ensures that the scaling, rotation, and translation adjustments are both precise and semantically informed.

*Step 3: Refinement.* With the updated partial point cloud  $p^{(k+1)}$ , the system could estimate a new similarity transformation to refine the alignment of the generated geometry within the scene. This updated partial point cloud is then fed back into the **Object Generation** module for the next iteration, allowing for progressive enhancements in both geometry accuracy and spatial positioning.

This iterative loop—alternating between geometry generation and transformation estimation—continues until convergence criteria are met. Convergence is achieved when the changes in transformation parameters fall below a predefined threshold or when a maximum number of iterations is reached. The result is a high-fidelity 3D object that is both visually accurate and geometrically aligned with the input data. By tightly integrating the **Object Generation** and **Alignment Generation** modules within an iterative framework, our approach effectively balances aesthetic fidelity with geometric precision. This joint generation process leverages both visual and depth information, ensuring that each 3D asset is of high quality and accurately positioned. Consequently, the pipeline lays a robust foundation for constructing physically correct and visually coherent 3D scenes, facilitating a wide range of downstream applications such as editing, rendering, and animation.

Once the object geometry is determined, we apply a state-of-the-art texture generation module to create photo-realistic surface details. Following established texture synthesis pipelines [Zhang et al. 2023, 2024a], we assign UV mappings and train a generative network to paint detailed textures onto the 3D meshes. This module is designed to robustly handle images under various augmentation, ensuring that the final textures match the input appearance even under occlusion or limited visibility.

## 5 PHYSICS-AWARE CORRECTION

The pipeline detailed in Sec. 4 individually generates each 3D object instance and estimates its similarity transformation (scaling, rotation, and translation) based on a single input image. While our proposed modules achieve high accuracy, the resulting scenes are sometimes not physically plausible. For instance, as illustrated in Fig. 4, one object (e.g., a guitar) may intersect with another (e.g., a cooler), or an object (e.g., a surfboard) may appear to float unnaturally without any support (e.g., from a van).

To address these issues, we introduce a physics-aware correction process that optimizes the rotation and translation of objects, ensuring the scene adheres to physical constraints consistent with common sense. The correction process is motivated by physical simulation (Sec. 5.1) and formulated as an optimization problem (Sec. 5.2) based on inter-object relations represented by a scene graph (Sec. 5.3) extracted from the image.

### 5.1 A Quick Primer to Rigid-Body Simulation

We introduce the fundamental principles of physical (rigid-body) simulation, which inspire our problem formulation and make our framework more accessible to downstream applications such as gaming and robotics. For a thorough survey, we refer the readers to Bender et al. [2012].

In rigid-body simulations, the world is modeled as an ordinary differential equation (ODE) process. In each simulation step, it begins with the Newton-Euler (differential) equations, which describe the dynamic motion of rigid bodies in the absence of contact. *Collision detection* is conducted to find the contact points between rigid bodies, which are needed to determine contact forces. For contact handling and collision resolution, there are usually several conditions: non-penetration constraints to prevent bodies from overlapping, a friction model ensuring contact forces remain within their friction cones, and complementarity constraints that enforce specific disjunctive relations among variables. Solvers are used to resolve the system comprising equations and inequalities, subsequently updating the velocity and position of each rigid body.

A straightforward approach to enhance physical plausibility is to utilize an off-the-shelf rigid-body simulator to process the scene, starting from the initial state estimated by the pipeline previously described and obtaining the rest state after simulation. However, this method presents several challenges.

1) *Partial Scene:* Some objects may be missing due to the limitations of 2D foundation models, and thus not reconstructed. Simulating a partial scene under full physical rules can lead to suboptimal results (see Fig. 10).

2) *Imperfect Geometries:* While our 3D generative model produces high-quality geometries, minor imperfections may still occur. Rigid-body simulators typically require convex decomposition [Mamou and Ghorbel 2009; Mamou et al. 2016; Wei et al. 2022] of objects, which introduces additional complexity and hyperparameters. Overly fine-grained decomposition can result in non-flat, complex surfaces, causing objects to fall or move unexpectedly during simulation. Conversely, coarse decomposition may lead to visually floating objects due to discrepancies between the visual and collision geometries.

3) *Initial Penetrations:* Despite the high accuracy of pose estimation, significant inter-object penetrations may exist in the initial state. These penetrations create instability for standard rigid-body solvers and, in some cases, lead to unsolvable scenarios if the solver is not customized for those cases.

Thus, we propose a customized and simplified “physical simulation” to optimize the object poses, ensuring that the scene adheres to common-sense physical principles derived from the single image. Note that our approach does not model full dynamics. For example, an object may not remain stable in its current pose over time. However,

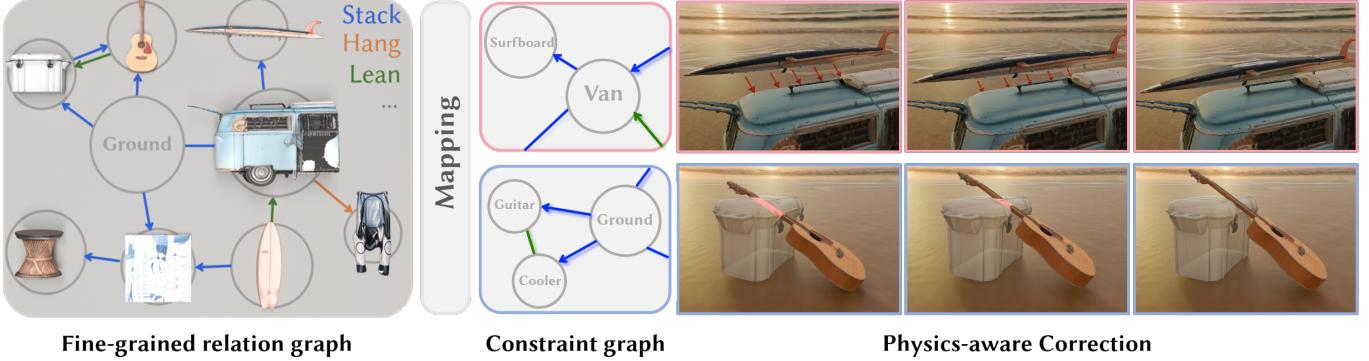


Fig. 4. Physics-aware correction via constraint graph mapped from fine-grained relation graph. Top: Floating surfboard grounded on the van. Bottom: Penetrating guitar and cooler separated.

it should be physically plausible at the current time step. We argue that our optimized results can serve as a reliable initialization for subsequent physical simulations.

## 5.2 Problem Formulation and Physical Constraints

We formulate the physics-aware correction process as an optimization problem, aiming to minimize the total cost that represents pairwise constraints on objects.

$$\min_{T=\{T_1, T_2, \dots, T_N\}} \sum_{i,j} C(T_i, T_j; o_i, o_j) \quad (8)$$

where  $N$  is the number of objects,  $T_i$  is the rigid transformation (rotation and translation) of the  $i$ -th object  $o_i$ .  $C$  is the cost function representing the relation between  $o_i$  and  $o_j$ . Note that the cost function varies depending on the type of relation.

Motivated by physical simulation, we categorize the relations into two types: *contact* and *support*. The relations are identified with the assistance of a VLM, as detailed in Sec. 5.3.

1) *Contact* describes whether two objects  $o_i$  and  $o_j$  are in contact. Let  $D_i(p)$  denote the signed distance function induced by  $o_i$  at the point  $p$ , which is used to define the constraint.  $D_i(p) = D_j(p) = 0$  indicates that  $p$  is a contact point of  $o_i$  and  $o_j$ . When  $D_i(p) = 0$  ( $p$  is a surface point of  $o_i$ ),  $D_j(p) < 0$  indicates inter-object penetration, while  $D_j(p) > 0$  means the objects are separated. Thus, the cost function can be defined as:

$$\begin{aligned} C(T_i, T_j; o_i \rightarrow o_j) &= -\frac{\sum_{p \in \partial o_j} D_i(p(T_j)) \mathbb{I}(D_i(p(T_j)) < 0)}{\sum_{p \in \partial o_j} \mathbb{I}(D_i(p(T_j)) < 0)} \\ &\quad + \max(\min_{p \in \partial o_j} D_i(p(T_j)), 0) \\ C(T_i, T_j; o_j \rightarrow o_i) &= -\frac{\sum_{p \in \partial o_i} D_j(p(T_i)) \mathbb{I}(D_j(p(T_i)) < 0)}{\sum_{p \in \partial o_i} \mathbb{I}(D_j(p(T_i)) < 0)} \\ &\quad + \max(\min_{p \in \partial o_i} D_j(p(T_i)), 0) \\ C(T_i, T_j) &= C(T_i, T_j; o_i \rightarrow o_j) + C(T_i, T_j; o_j \rightarrow o_i) \\ &\quad \text{if } o_i \text{ and } o_j \text{ are in contact} \end{aligned} \quad (9)$$

where  $\partial o_i$  denotes the surface of  $o_i$ , and  $\mathbb{I}$  is the indicator function. The constraint ensures that there is no penetration and at least one

contact point between the objects. Note that  $p \in \partial o_i$  is a function of  $T_i$ . The contact constraint defined here is bilateral, meaning it applies to both objects.

2) *Support* is a unilateral constraint, which is a special case of *Contact*. If  $o_i$  supports  $o_j$ , it implies that the pose  $T_j$  of  $o_j$  should be optimized while  $o_i$  is assumed to be static. This scenario typically occurs when multiple objects are stacked vertically. The cost function for this case is similar to the one in *Contact*, but it only involves one direction:

$$C(T_i, T_j) = |\min_{p \in \partial o_j} D_i(p(T_j))|, \text{ if } o_i \text{ supports } o_j \quad (10)$$

Furthermore, for flat supporting surfaces like the ground or walls, we regularize the SDF values near the contact region, to ensure that objects make close contact with these surfaces. This regularization handles scenarios where objects are partially reconstructed, such as a van with only two wheels, as illustrated in Fig. 4.

$$C(T_i, T_j) = \frac{\sum_{p \in \partial o_j} D_i(p(T_j)) \mathbb{I}(0 < D_i(p) < \sigma)}{\sum_{p \in \partial o_j} \mathbb{I}(0 < D_i(p) < \sigma)} \quad (11)$$

where  $\mathbb{I}$  is the indicator function, and  $\sigma$  is a threshold to decide whether a point is sufficiently close to the surface.

## 5.3 Scene Relation Graph

Physical cues, particularly inter-object relations, are visually present in the image. We leverage the strong common-sense reasoning capabilities [Cheng et al. 2024; Li et al. 2024b; Rana et al. 2023] of visual-language models, specifically GPT-4v [Achiam et al. 2023], to identify pairwise physical constraints as defined in Sec. 5.2. Given an image, we employ the Set of Mark [Yang et al. 2023] (SoM) technique to visually prompt GPT-4v to describe the inter-object relations, and subsequently extract a *scene relation graph* from the answers. To address the sampling uncertainty inherent in VLMs, we adopt an ensemble strategy, combining results from multiple trials. We define relations as correct if they appear in more than half of the samples to produce a robust inferred graph. To be more specific, we apply the Set-of-Mark method with random colorization and numerical ordering multiple times, enabling more reliable and consistent outputs for further GPT-based question-answering tasks.

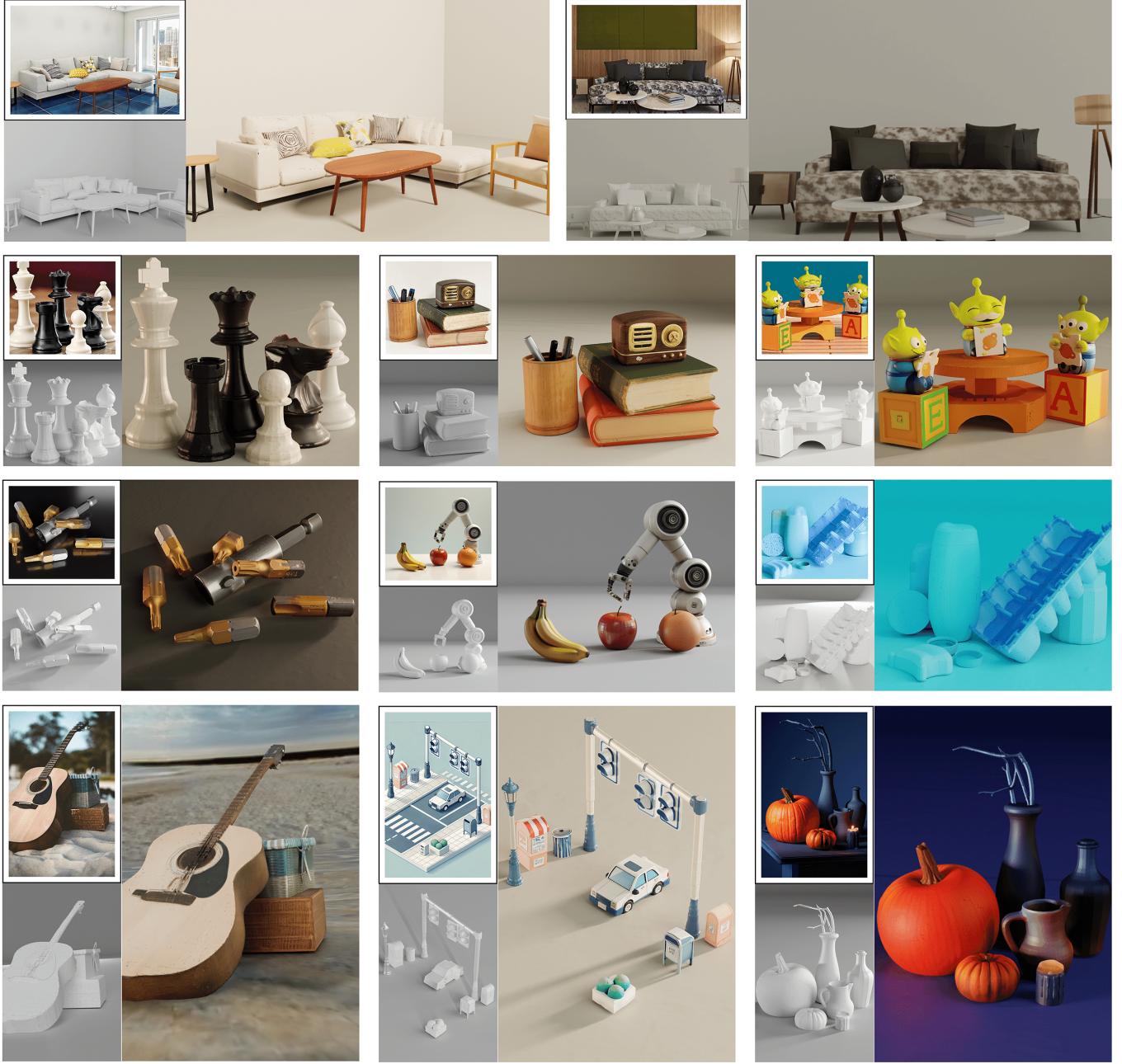


Fig. 5. Bringing the vibrant diversity of the real world into the virtual realm, this collection reimagines open-vocabulary scenes as immersive digital environments, capturing the richness and depth of each unique setting. For each scene, the images display as follows: the top-left shows the input image, the top-center displays the rendered geometry, and the right presents the rendered image with realistic textures.

Instead of directly asking GPT-4v to identify *Support* and *Contact* relations, we first provide it with more fine-grained physical relations, such as *Stack* (Object 2 supports Object 1), *Lean* (Object 1 leans against Object 2), and *Hang* (Object 2 supports Object 1 from above). We instruct GPT-4v to analyze numbered objects from the Set-of-Mark method and output all contact-based relations, covering

six types: Stack, Lean, Hang, Clamped, Contained, and Edge/Point. The prompt specifies that only contacting objects have relations and defaults to Stack for ambiguous cases.

We then map these detailed relations to the predefined categories of *Support* and *Contact* for further optimization. Specifically, if there are edges pointing toward each other between two nodes, the edge

is categorized as *Contact*; otherwise, it is categorized as *Support*. Prompting GPT-4v with these nuanced relations helps eliminate potential ambiguity in binary relation classification and facilitates more accurate reasoning by GPT-4v. An example of the resulting graph is illustrated in Fig. 4.

The mapped scene constraint graph is a directed graph where nodes represent object instances and edges denote physical relations between objects. A *Contact* relation is represented by a bidirectional edge, while a *Support* relation is depicted as a directed edge. This graph serves as the foundation for defining the cost functions used in Eq. 8.

#### 5.4 Optimization with Physics-Aware Relation Graph

Given the physical constraints defined by the inferred relation graph, we can instantiate our cost functions as described in Eq. 8. The graph allows us to reduce the number of pairwise constraints that need to be optimized, in contrast to a full physical simulation.

For the implementation, we uniformly sample a fixed number of points from the surface of each object at its rest pose. These points are then transformed according to the current object’s pose parameters and used to query the SDF values with respect to another object (and its pose). SDF computation is handled by Open3D, and Pytorch is used to auto-differentiate the loss function.

## 6 RESULT

Fig. 5 showcases a range of 3D scenes generated by our method from single-view inputs across a diverse range of open-vocabulary scenarios, featuring detailed indoor environments, close-up captures of objects, and AI-generated imagery. These examples highlight the versatility and robustness of our approach, exhibiting high-fidelity geometry, realistic textures, and convincing scene compositions.

### 6.1 Implementation Details

*ObjectGen.* The ObjectGen (Sec. 4.1) model’s pretraining follows the methodology outlined in 3DShape2VecSet [Zhang et al. 2023] and CLAY [Zhang et al. 2024a], where we leverage both a Variational Autoencoder (VAE) and a Latent Diffusion Model (LDM) to generate 3D object geometries. Both the VAE and LDM modules are implemented using a 24-layer transformer, comprising a total of 1.5 billion parameters. The model is trained on the Objaverse [Deitke et al. 2023] dataset, which consists of approximately 500,000 3D assets after filtering. The partial point cloud conditioning follows a similar approach to CLAY’s adaptation framework. We encode the canonical-space partial point cloud as positional embeddings with a feature dimension of 512, which are injected into the main LDM transformer using cross-attention mechanisms. For each 3D asset, we render 32 views and precompute depth maps using MoGe [Wang et al. 2024b] and Metric3D [Yin et al. 2023]. These depth maps are then lifted into point clouds during training, with random masks applied to simulate occlusions. We sample 2048 points from the unprojected point cloud using Farthest Point Sampling (FPS), which serves as conditioning input for the LDM. To enhance the model’s robustness, we randomly interpolate between the ground truth and predicted partial point clouds, allowing the system to handle data of varying quality. The conditioning module is trained on 200K curated

Table 1. Quantitative comparison of scene reconstruction methods across four metrics with CLIP score, GPT-4 ranking, user study of visual quality (VQ), and physical plausibility (PP).

Method	CLIP↑	GPT-4↓	VQ↑	PP↑
ACDC	69.77	2.7	5.58%	22.86%
Gen3DSR	79.84	2.175	6.35%	5.72%
ours	<b>85.77</b>	<b>1.125</b>	<b>88.07%</b>	<b>71.42%</b>

Table 2. Quantitative comparison of scene reconstruction performance on the 3D-Front indoor dataset. We evaluate different methods based on Chamfer Distance (CD) for shape accuracy, F-Score (FS) for object-level reconstruction quality, and Intersection over Union (IoU) for scene-level overlap.

Method	CD-S↓	FS-S↑	CD-O↓	FS-O↑	IoU-B↑
ACDC	0.104	39.46	0.072	41.99	0.541
InstPIFU	0.092	39.12	0.103	38.29	0.436
Gen3DSR	0.083	38.95	0.071	39.13	0.459
ours	<b>0.052</b>	<b>56.18</b>	<b>0.057</b>	<b>56.50</b>	<b>0.603</b>

data from Objaverse over 3000 epochs with 64 Nvidia A800 GPUs required approximately one week. The AdamW optimizer is used with a learning rate of 1e-5. For inferencing a single object, object generation takes approximately 7 seconds, and texture generation takes approximately 10 seconds per object on an NVIDIA A6000 GPU.

*AlignGen.* The AlignGen (Sec. 4.2) module, responsible for generating pose alignment, utilizes a 24-layer transformer with a feature dimension of 512, resulting in a total of 150 million parameters. During training, we randomly sample a partial point cloud in the canonical space from the point clouds lifted from precomputed depth maps and apply a random transformation to this point cloud. The transformed point cloud, along with the geometry latent code  $z$  from ObjectGen, is used as the conditioning input. FPS is utilized to sample 2048 points from the partial point cloud to ensure a fixed number of inputs to the transformer. Training is conducted on the same 200K curated dataset over 1,500 epochs with 64 Nvidia A800 GPUs for approximately two day. The AdamW optimizer is used with a learning rate of 1e-5. During inference, AlignGen module takes around 1 second for one object for pose generation.

### 6.2 Comparison

*Qualitative Comparisons.* We first evaluate our method, CAST, against state-of-the-art single-image scene reconstruction techniques on open-vocabulary scenarios. We also included images used by ACDC and Gen3DSR to further demonstrate the scene reconstruction results of different methods. Fig. 6 illustrates the performance of three methods—(1) the retrieval-based approach ACDC [Dai et al. 2024], (2) the generation-based method Gen3DSR [Dogaru et al. 2024], and (3) our proposed CAST—across both reference and novel views. Our results highlight CAST’s superior ability to accurately reconstruct scenes in diverse settings, including indoor and outdoor environments, close-up perspectives, and AI-generated imagery.

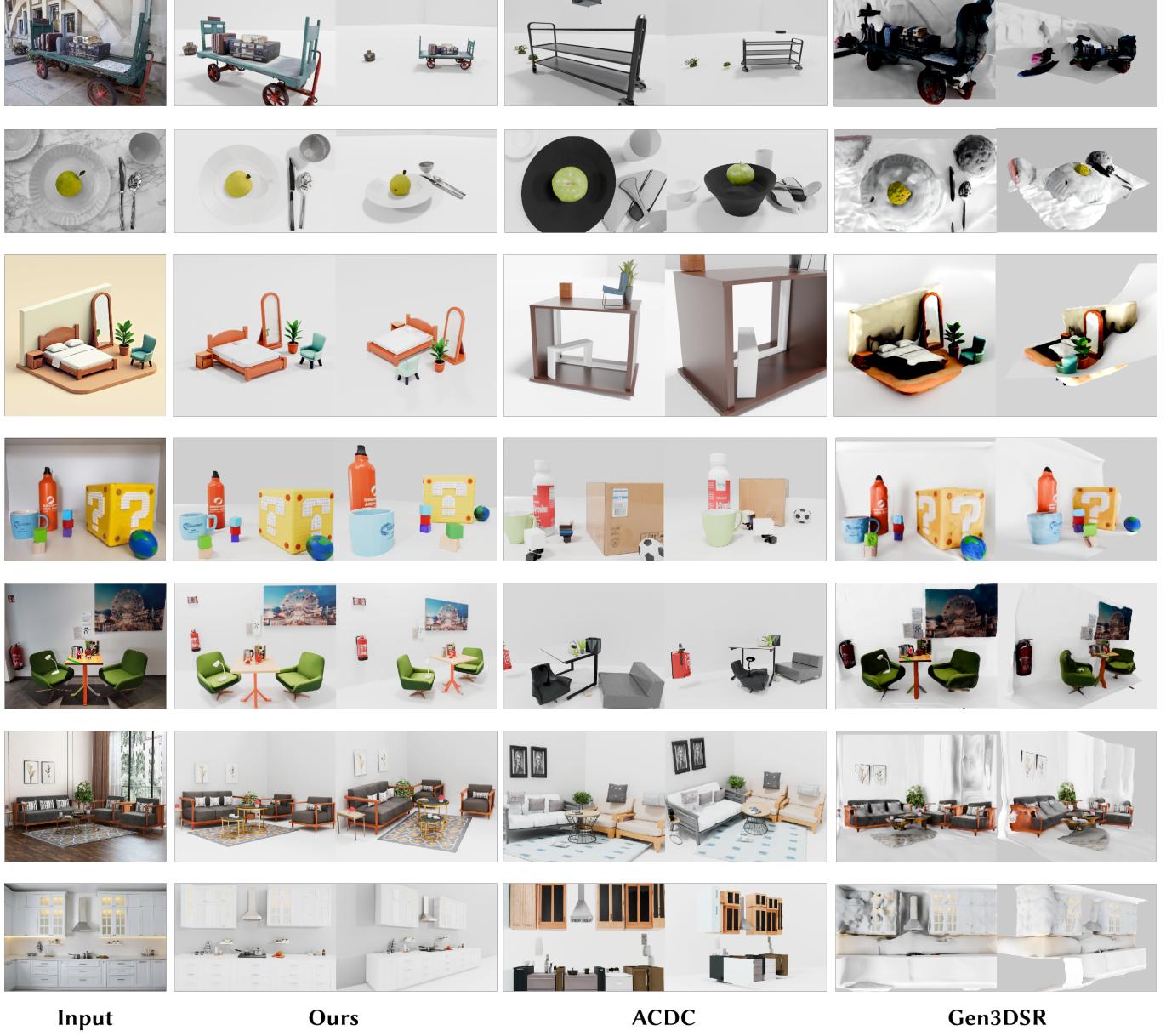


Fig. 6. Qualitative comparisons of CAST with state-of-the-art single-image scene reconstruction methods. From left to right: Input image, CAST, ACDC, and Gen3DSR. Top to bottom: random open vocabulary dataset (rows 1–3), Gen3DSR input (rows 4–5), ACDC input (rows 6–7).

As shown in Fig.6, CAST distinguishes itself from both ACDC and Gen3DSR through innovative advancements. Unlike ACDC, which is limited to indoor scenes and relies on large datasets for object retrieval, often producing objects similar to those in the scene rather than the objects themselves, CAST supports open-vocabulary generalization. This allows CAST to accurately reconstruct objects in varied and complex environments. While ACDC uses simple bounding boxes as proxies, CAST combines image-based physical priors with mesh optimization to effectively manage complex scenes. In contrast to Gen3DSR, CAST employs direct 3D generation via a

Masked Autoencoder, eliminating the error-prone 2D inpainting step. This results in smoother meshes, significantly outperforming Gen3DSR in single-object generation quality, particularly in challenging scenes. Moreover, Gen3DSR's lack of simulation often leads to issues such as interpenetration or floating objects, causing scenes to appear consistent only from the input viewpoint and degrading novel view rendering. CAST, by contrast, ensures robust scene consistency across perspectives. CAST demonstrates robust scene reconstructions under varying conditions, underscoring its versatility for a broad range of real-world and generated scenarios.

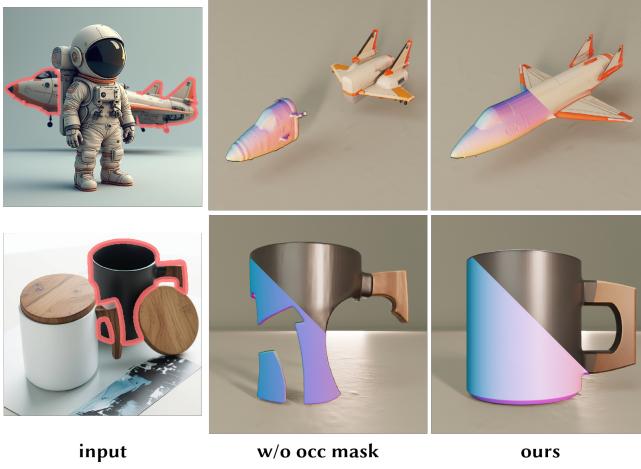


Fig. 7. We evaluate the generation performance with and without the occlusion-aware generation module. The RGB and normal renderings of the object highlight the significance of this module in ensuring the completeness and high quality of the generated object.

To assess both the visual fidelity and semantic accuracy of the generated scenes, we employ two complementary evaluation methods including CLIP Score [Zhengwentai 2023] and GPT-4 Reasoning. We compute the CLIP score between the rendered scene and the input image to measure overall reconstruction quality and visual similarity. To minimize environmental distractions, we remove backgrounds from both the rendered and reference images before computing the score. We additionally leverage GPT-4 to rank the generated scenes based on various semantic aspects, including object arrangement, physical relations, and scene realism. This semantic feedback helps identify alignment or contextual errors that might not be apparent through pixel-based scores alone.

Beyond automated metrics, we also conducted a user study focusing on two key aspects including Visual Quality (VQ) and Physical Plausibility (PP). We randomly selected paired reference, novel, and target views, asking participants to choose which method’s output best matched the input image in terms of both similarity and overall aesthetics. To reduce potential biases introduced by visual resemblance, participants in a separate session only viewed rendered results—without the original input images—and judged which scene appeared more realistic based on physical constraints and common sense (e.g., preventing floating objects or improbable contacts).

As shown in Tab. 1, CAST outperforms both ACDC and Gen3DSR in all four evaluated metrics, confirming its effectiveness in producing scenes that are visually coherent and physically plausible.

*Quantitative Comparisons.* Although CAST is designed to handle open-vocabulary scenes, many such scenes lack mesh ground truth, which complicates direct quantitative comparisons. To address this, we perform additional evaluations on the 3DFront dataset [Fu et al. 2021]. This dataset offers ground-truth meshes alongside corresponding rendered images, enabling a more precise assessment of both object-level and scene-level reconstructions. We compare our method with InstPIFu [Liu et al. 2022], ACDC [Dai et al. 2024], and

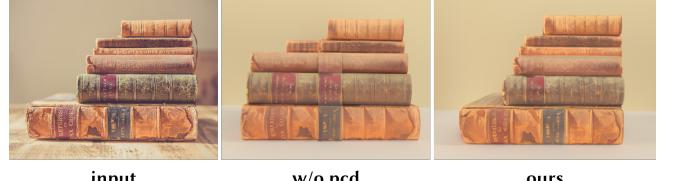


Fig. 8. A stack of books with varying lengths and widths directly generated as a single complex object, demonstrating how point cloud conditioning enhances the preservation of scale, dimensions, and local details compared to traditional methods.

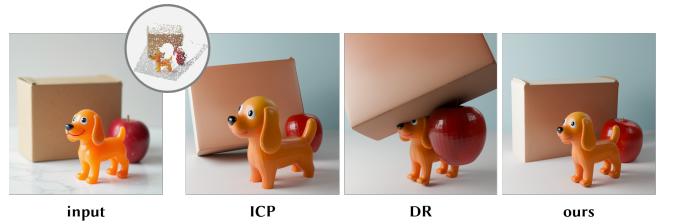


Fig. 9. Comparative evaluation of pose estimation methods. Our pose alignment module demonstrates superior alignment accuracy compared to Iterative Closest Point (ICP) and differentiable rendering (DR).

Gen3DSR [Dogaru et al. 2024]. We compute Chamfer Distance and F-Score at the object level, as well as IoU, Chamfer Distance, and F-Score at the scene level, to assess both the fidelity of individual object geometries and the accuracy of their spatial layout. To ensure fairness, we replace the segmentation modules in other methods with ground-truth (GT) masks so that any differences stem purely from reconstruction ability rather than object partitioning.

As summarized in Tab. 2, CAST not only achieves higher object-level generation quality but also surpasses existing approaches in scene layout accuracy. Even within the constraints of an indoor dataset, our method demonstrates robust performance and consistent improvements over competing baselines.

### 6.3 Evaluations

To elucidate the individual contributions of key components in CAST, we conducted a series of ablation studies. These experiments systematically removed or altered specific components to assess their impact on overall performance. The ablation studies focus on several key design choices: occlusion-aware object generation, point-cloud conditioning, generative alignment, and the physics-aware correction process.

*Ablation on Occlusion-Aware Generation.* Occlusions are a significant challenge in scenes with complex objects. To assess the effectiveness of the Masked Autoencoder (MAE) in handling occlusions, we performed an ablation study comparing generation results with and without MAE components. As shown in Fig. 7, the results highlight the importance of the occlusion-aware module. Without MAE, the generated objects for partially occluded regions exhibit significant degradation. For example, the spaceship appears fragmented and incomplete, while the cup is depicted as broken

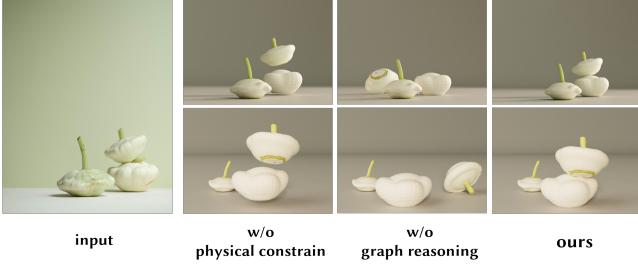


Fig. 10. Comparison of scene reconstruction with and without relational graph constraints. By integrating relational graph constraints, our method ensures both physical plausibility and accurate alignment with the intended scene, maintaining correct spatial relations.

with missing parts. In contrast, when MAE conditioning is applied, the model successfully infers and fills the occluded regions, resulting in more accurate and visually coherent generations that align better with the input image. This demonstrates the critical role of the occlusion-aware module in ensuring that occluded objects are reconstructed accurately, improving both the completeness and realism of the final 3D scene.

*Ablation on Partial Point Cloud Conditioning.* We conducted an ablation study to investigate the role of canonical space partial point cloud conditioning in our generation process. Although directly generating from the input image can produce visually plausible results. In the absence of pixel-level alignment, the model struggles to maintain correct object quantity and scale, resulting in unsatisfactory generation. To more effectively showcase the importance of point cloud conditioning in generating a single instance, we opted to directly generate a more complex instance structure: a stack of six books with varying lengths and widths. As depicted in Fig. 8. When the generation process relies solely on the input image, without the benefit of point cloud conditioning, the results frequently exhibit inaccuracies in both the number and dimensions of the generated objects. By contrast, integrating point cloud conditioning introduces a robust geometric prior that significantly improves the precision of the generated scene. This enhancement ensures that objects with intricate shapes and varying dimensions are reconstructed more accurately, closely resembling their real-world counterparts depicted in the input image. This demonstrates the critical role of geometric priors in enhancing the fidelity of 3D scene generation by preserving the true dimensions and shapes.

*Effectiveness of Alignment Generation.* To assess the effectiveness of our pose alignment module, we compared it with common pose estimation methods such as Iterative Closest Point (ICP) [Arun et al. 1987; Best 1992] and differentiable rendering [Laine et al. 2020]. The generated mesh was provided to different pose estimation methods to align it with the reference RGB image and its corresponding depth prediction. For the ICP method, we uniformly sampled a point cloud from the generated mesh and normalized both the sampled and estimated point clouds by their bounding boxes to address scale differences. We used the ICP implementation in Open3D [Zhou

Table 3. Quantitative ablation study of the MAE module, point cloud conditioning (PCD), and the iterative refinement strategy (iter.). For simplicity, we only display the added key component in each row.

Method	CD-S↓	FS-S↑	CD-O↓	FS-O↑	IoU-B↑
Vanilla	0.079	53.38	0.069	52.83	0.515
+ MAE	0.064	53.79	0.066	54.32	0.548
+ PCD	0.056	53.91	0.060	54.60	0.582
+ iter.	<b>0.052</b>	<b>56.18</b>	<b>0.057</b>	<b>56.50</b>	<b>0.603</b>

et al. 2018] to register these two normalized point clouds. For differentiable rendering, we optimized the rotation and translation parameters to transform the generated mesh so that the rendered image aligned with the reference RGB image. As shown in Fig. 9, our method surpasses both ICP and differentiable rendering in alignment accuracy. ICP often struggles with accurate pose estimation due to outliers in point clouds, unknown object scales, and symmetrical or repetitive geometries, which can lead to local minima. Differentiable rendering, on the other hand, is significantly impacted by occlusions in the RGB input, disrupting the optimization of object poses and preventing precise alignment with the input image. Our results show that our pose alignment module outperforms traditional ICP and differentiable rendering methods, demonstrating its robustness in accurately estimating object poses from generated meshes and improving alignment with input images.

*Effect of Physical Consistency Enforcement.* In CAST, physical constraints are essential for achieving realistic object interactions and maintaining spatial coherence within a scene. While we address common challenges such as occlusions and incomplete views, issues like floating objects, penetration, and misaligned spatial relations still occur. As shown in Fig. 10, scenes generated without relational constraints may appear physically inconsistent, when only physical simulation is applied, objects adhere to physical laws, but their relative positioning and overall arrangement can differ significantly from the intended scene (e.g., an onion might fall off a surface, disrupting the original composition). By incorporating relational graph constraints, our method ensures that the objects not only comply with physical feasibility but also align with the intended scene layout, preserving both the physical plausibility and the desired spatial relations.

*Quantitative Ablation Study of Different Modules.* To quantitatively evaluate the contribution of each module to overall performance, we conducted a comprehensive ablation study. As shown in Tab. 3, we assessed the impact of removing or altering key components on the final scene quality. The results indicate that each component contributes significantly to the overall performance of our method. The quantitative analysis further highlights the importance of each module in achieving high-quality, physically consistent, and realistic scene reconstructions.

*Applications.* As shown in Fig. 11, CAST transforms a single image into a fully realized 3D scene, enabling a wide range of applications. This ability to reconstruct detailed environments powers physics-based animation by ensuring realistic object interactions. It



Fig. 11. CAST enables realistic physics-based animations, immersive game environments, and efficient real-to-simulation transitions, driving innovation across various fields.

also supports real-to-simulation workflows in robotics, allowing for accurate scene replication from real-world datasets. In game development, CAST facilitates the creation of immersive environments, where faithfully reconstructed scenes are seamlessly integrated into interactive worlds using Unreal Engine.

## 7 CONCLUSIONS

In this paper, we introduced CAST, a novel single-image 3D scene reconstruction method that combines geometric fidelity, pixel-level alignment, and physically grounded constraints. By integrating scene decomposition, a perceptive 3D instance generation framework, and physical correction techniques, CAST addresses key challenges such as pose misalignment, object interdependencies, and partial occlusions. This structured pipeline results in 3D scenes that are both visually accurate and physically consistent, pushing beyond the limitations of traditional object-centric approaches. We validated CAST through extensive experiments and user studies, demonstrating significant performance improvements over state-of-the-art methods in terms of visual quality and physical plausibility. We anticipate that CAST will serve as a strong foundation for future developments in 3D generation, scene reconstruction, and immersive content creation.

*Limitations and Future Work.* The quality of scene generation in CAST is heavily dependent on the underlying object generation model. At present, the model still lacks sufficient detail and precision, this limitation leads to noticeable inconsistencies in the generated objects, affecting their alignment and spatial relations in the scene.



Fig. 12. In some scenes, transparent glass, textiles, and fabrics are difficult to express, as the mesh struggles to represent them realistically.

Additionally, the current mesh representation struggles with materials like textiles, glasses or fabrics, often appearing unnatural, and fails to accurately depict transparent materials, as shown in Fig. 12. Although additional modules have been incorporated to enhance object robustness and similarity, the need for more advanced and robust generation models remains. A more detailed and accurate object generator could significantly improve the overall scene quality and enhance its real-world applicability.

A notable limitation of the current method is the absence of lighting estimation and background modeling. Without realistic lighting, the interactions between objects and their surroundings may lack natural shading and illumination effects, impacting the visual realism and immersion of the generated 3D environments. For enhanced visual realism, we employ an off-the-shelf panoramic HDR generation tool [Hyper3D 2025] combined with preset lighting conditions in blender manually. Future enhancements in CAST could benefit from integrating advanced techniques for lighting estimation and background modeling, which would significantly enrich the contextual depth and visual fidelity of the scenes.

In more complex scenes, the performance of the current method may experience slight degradation. Challenges such as intricate spatial layouts and dense object configurations could affect the accuracy of scene reconstruction to some extent. While CAST currently excels at reconstructing individual scenes, there is significant potential to utilize its outputs to build large-scale datasets, facilitating advanced research on fully learned scene or video generation pipelines. Expanding the variety and realism of generated scenes in this manner could further improve the robustness and applicability of 3D generative models in areas such as film production, simulation, and immersive media.

## ACKNOWLEDGMENTS

This work was supported by National Key R&D Program of China (2022YFF0902301), NSFC programs (61976138, 61977047), STCSM (2015F0203-000-06), and SHMEC (2019-01-07-00-01-E00003). We also acknowledge support from Shanghai Frontiers Science Center of Human-centered Artificial Intelligence (ShangHAI), MoE Key Lab of Intelligent Perception and Human-Machine Collaboration (ShanghaiTech University), Core Facility Platform of Computer Science and Communication of ShanghaiTech University, and HPC Platform of ShanghaiTech University.

## REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).
- K Somani Arun, Thomas S Huang, and Steven D Blostein. 1987. Least-squares fitting of two 3-D point sets. *IEEE Transactions on pattern analysis and machine intelligence* 5 (1987), 698–700.
- Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. 2021. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF international conference on computer vision*. 5855–5864.
- Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. 2022. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5470–5479.
- Jan Bender, Kenny Erleben, Jeff Trinkle, and Erwin Coumans. 2012. Interactive Simulation of Rigid Body Dynamics in Computer Graphics. In *33rd Annual Conference of the European Association for Computer Graphics, Eurographics 2012 - State of the Art Reports, Cagliari, Sardinia, Italy, May 13-18, 2012*, Marie-Paule Cani and Fabio Ganovelli (Eds.). Eurographics Association, 95–134. <https://doi.org/10.2312/CONF-EG2012-STARS/095-134>
- Paul J Best. 1992. A method for registration of 3-D shapes. *IEEE Trans Pattern Anal Mach Vision* 14 (1992), 239–256.
- Shariq Farooq Bhat, Reiner Birk, Diana Wofk, Peter Wonka, and Matthias Müller. 2023. Zoedepth: Zero-shot transfer by combining relative and metric depth. *arXiv preprint arXiv:2302.12288* (2023).
- Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. 2023. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127* (2023).
- Jake Bruce, Michael D Dennis, Ashley Edwards, Jack Parker-Holder, Yuge Shi, Edward Hughes, Matthew Lai, Aditi Mavalankar, Richie Steigerwald, Chris Apps, et al. 2024. Genie: Generative interactive environments. In *Forty-first International Conference on Machine Learning*.
- Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. 2015. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012* (2015).
- Anpei Chen, Minye Wu, Yingliang Zhang, Nianyi Li, Jie Lu, Shenghua Gao, and Jingyu Yu. 2018. Deep surface light fields. *Proceedings of the ACM on Computer Graphics and Interactive Techniques* 1, 1 (2018), 1–17.
- Yixin Chen, Junfeng Ni, Nan Jiang, Yaowei Zhang, Yixin Zhu, and Siyuan Huang. 2024a. Single-view 3d scene reconstruction with high-fidelity shape and texture. In *2024 International Conference on 3D Vision (3DV)*. IEEE, 1456–1467.
- Yunuo Chen, Tianyi Xie, Zeshun Zong, Xuan Li, Feng Gao, Yin Yang, Ying Nian Wu, and Chenfanfu Jiang. 2024b. Atlas3D: Physically Constrained Self-Supporting Text-to-3D for Simulation and Fabrication. *arXiv preprint arXiv:2405.18515* (2024).
- An-Chieh Cheng, Yandong Ji, Zhaojing Yang, Xueyan Zou, Jan Kautz, Erdem Biyik, Hongxu Yin, Sifei Liu, and Xiaolong Wang. 2024. Navila: Legged robot vision-language-action model for navigation. *arXiv preprint arXiv:2412.04453* (2024).
- Tao Chu, Pan Zhang, Qiong Liu, and Jiaqi Wang. 2023. Buol: A bottom-up framework with occupancy-aware lifting for panoptic 3d scene reconstruction from a single image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4937–4946.
- Manuel Dahmert, Ji Hou, Matthias Nießner, and Angela Dai. 2021. Panoptic 3d scene reconstruction from a single rgb image. *Advances in Neural Information Processing Systems* 34 (2021), 8282–8293.
- Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. 2017. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 5828–5839.
- 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 Laforte, Vikram Voleti, Samir Yitzhak Gadre, et al. 2024. Objaverse-xl: A universe of 10m+ 3d objects. *Advances in Neural Information Processing Systems* 36 (2024).
- Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. 2023. Objaverse: A universe of annotated 3d objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 13142–13153.
- Andreea Dogaru, Mert Özer, and Bernhard Egger. 2024. Generalizable 3D Scene Reconstruction via Divide and Conquer from a Single View. *arXiv preprint arXiv:2404.03421* (2024).
- Huan Fu, Bowen Cai, Lin Gao, Ling-Xiao Zhang, Jiaming Wang, Cao Li, Qixun Zeng, Chengyue Sun, Rongfei Jia, Binqiang Zhao, et al. 2021. 3d-front: 3d furnished rooms with layouts and semantics. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 10933–10942.
- Daoyi Gao, Dávid Rozemberczki, Stefan Leutenegger, and Angela Dai. 2024b. Diffcad: Weakly-supervised probabilistic cad model retrieval and alignment from an rgb image. *ACM Transactions on Graphics (TOG)* 43, 4 (2024), 1–15.
- Ruiqi Gao, Aleksander Holynski, Philipp Henzler, Arthur Brussee, Ricardo Martin-Brualla, Pratul Srinivasan, Jonathan T Barron, and Ben Poole. 2024a. Cat3d: Create anything in 3d with multi-view diffusion models. *arXiv preprint arXiv:2405.10314* (2024).
- Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. 2013. Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research* 32, 11 (2013), 1231–1237.
- Georgia Gkioxari, Nikhila Ravi, and Justin Johnson. 2022. Learning 3d object shape and layout without 3d supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1695–1704.
- Michael Goesele, Noah Snavely, Brian Curless, Hugues Hoppe, and Steven M. Seitz. 2007. Multi-View Stereo for Community Photo Collections. In *2007 IEEE 11th International Conference on Computer Vision*. 1–8. <https://doi.org/10.1109/ICCV.2007.4408933>
- Can Gümel, Angela Dai, and Matthias Nießner. 2022. Roca: Robust cad model retrieval and alignment from a single image. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 4022–4031.
- Minghao Guo, Bohan Wang, Pingchuan Ma, Tianyuan Zhang, Crystal Elaine Owens, Chuang Gan, Joshua B Tenenbaum, Kaiming He, and Wojciech Matusik. 2024. Physically Compatible 3D Object Modeling from a Single Image. *arXiv preprint arXiv:2405.20510* (2024).
- Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. 2022a. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303* (2022).
- Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. 2022b. Video diffusion models. *Advances in Neural Information Processing Systems* 35 (2022), 8633–8646.
- Yicong Hong, Kai Zhang, Jieyang Gu, Sai Bi, Yang Zhou, Difan Liu, Feng Liu, Kalyan Sunkavalli, Trung Bui, and Hao Tan. 2023. Lrm: Large reconstruction model for single image to 3d. *arXiv preprint arXiv:2311.04400* (2023).
- Zehuan Huang, Yuan-Chen Guo, Xingqiao An, Yunhan Yang, Yangguang Li, Zi-Xin Zou, Ding Liang, Xihui Liu, Yan-Pei Cao, and Lu Sheng. 2024. MIDI: Multi-Instance Diffusion for Single Image to 3D Scene Generation. *arXiv preprint arXiv:2412.03558* (2024).
- Hyper3D. 2025. *Omnircraft*. <https://hyper3d.ai/omnircraft/hdr/>
- Wadim Kehl, Fabian Manhardt, Federico Tombari, Slobodan Ilic, and Nassir Navab. 2017. Ssd-6d: Making rgb-based 3d detection and 6d pose estimation great again. In *Proceedings of the IEEE international conference on computer vision*. 1521–1529.
- Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 2023. 3D Gaussian Splatting for Real-Time Radiance Field Rendering. *ACM Transactions on Graphics* 42, 4 (July 2023). <https://repo-sam.inria.fr/fungraph/3d-gaussian-splatting/>
- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. 2023. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 4015–4026.
- Weicheng Kuo, Anelia Angelova, Tsung-Yi Lin, and Angela Dai. 2021. Patch2cad: Patchwise embedding learning for in-the-wild shape retrieval from a single image. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 12589–12599.
- Yann Labb  , Justin Carpentier, Mathieu Aubry, and Josef Sivic. 2020. Cosopose: Consistent multi-view multi-object 6d pose estimation. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XVII*. Springer, 574–591.
- Samuli Laine, Janne Hellsten, Tero Karras, Yeongho Seol, Jaakko Lehtinen, and Timo Aila. 2020. Modular primitives for high-performance differentiable rendering. *ACM Transactions on Graphics (TOG)* 39 (2020), 1 – 14.

- Florian Langer, Gwangbin Bae, Ignas Budvytis, and Roberto Cipolla. 2022. SPARC: Sparse render-and-compare for CAD model alignment in a single RGB image. *arXiv preprint arXiv:2210.01044* (2022).
- Bruno Latour. 2005. *Reassembling the Social: An Introduction to Actor-Network-Theory*. Oxford University Press, Oxford, UK.
- Wenhao Li, Zhiyuan Yu, Qijin She, Zhihan Yu, Yuqing Lan, Chenyang Zhu, Ruizhen Hu, and Kai Xu. 2024b. LLM-enhanced Scene Graph Learning for Household Rearrangement. In *SIGGRAPH Asia 2024 Conference Papers*. 1–11.
- Xuanlin Li, Kyle Hsu, Jiayuan Gu, Karl Pertsch, Oier Mees, Homer Rich Walke, Chuyuan Fu, Ishikaa Lunawat, Isabel Sieh, Sean Kirmani, et al. 2024a. Evaluating Real-World Robot Manipulation in Simulation. *arXiv preprint arXiv:2405.05941* (2024).
- Yixun Liang, Xin Yang, Jiantao Lin, Haodong Li, Xiaogang Xu, and Yingcong Chen. 2024. Luciddreamer: Towards high-fidelity text-to-3d generation via interval score matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6517–6526.
- Haolin Liu, Yujian Zheng, Guanying Chen, Shuguang Cui, and Xiaoguang Han. 2022. Towards high-fidelity single-view holistic reconstruction of indoor scenes. In *European Conference on Computer Vision*. Springer, 429–446.
- Minghua Liu, Chao Xu, Haian Jin, Linghao Chen, Mukund Varma T, Zexiang Xu, and Hao Su. 2024. One-2-3-45: Any single image to 3d mesh in 45 seconds without per-shape optimization. *Advances in Neural Information Processing Systems* 36 (2024).
- Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick. 2023c. Zero-1-to-3: Zero-shot one image to 3d object. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*.
- Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Qing Jiang, Chunyuan Li, Jianwei Yang, Hang Su, et al. 2025. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. In *European Conference on Computer Vision*. Springer, 38–55.
- Xueyi Liu, Bin Wang, He Wang, and Li Yi. 2023b. Few-Shot Physically-Aware Articulated Mesh Generation via Hierarchical Deformation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 854–864.
- Yuan Liu, Cheng Lin, Zijiao Zeng, Xiaoxiao Long, Lingjie Liu, Taku Komura, and Wenping Wang. 2023a. SyncDreamer: Generating Multiview-consistent Images from a Single-view Image. In *arXiv preprint arXiv:2309.03453*.
- Xiaoxiao Long, Yuan-Chen Guo, Cheng Lin, Yuan Liu, Zhiyang Dou, Lingjie Liu, Yuexin Ma, Song-Hai Zhang, Marc Habermann, Christian Theobalt, et al. 2024. Wonder3d: Single image to 3d using cross-domain diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9970–9980.
- Khaled Mamou and Faouzi Ghorbel. 2009. A simple and efficient approach for 3D mesh approximate convex decomposition. In *2009 16th IEEE international conference on image processing (ICIP)*. IEEE, 3501–3504.
- Khaled Mamou, E Lengyel, and A Peters. 2016. Volumetric hierarchical approximate convex decomposition. *Game engine gems 3* (2016), 141–158.
- Mariem Mezghanni, Théo Bodroito, Malika Boulkenafed, and Maks Ovsjanikov. 2022. Physical simulation layer for accurate 3d modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 13514–13523.
- Mariem Mezghanni, Malika Boulkenafed, Andre Lieutier, and Maks Ovsjanikov. 2021. Physically-aware generative network for 3d shape modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9330–9341.
- Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. 2020. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European conference on computer vision*. Springer, 405–421.
- Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. 2022. Instant Neural Graphics Primitives with a Multiresolution Hash Encoding. *ACM Trans. Graph.* 41, 4, Article 102 (July 2022), 15 pages. <https://doi.org/10.1145/3528223.3530127>
- Junfeng Ni, Yixin Chen, Bohan Jing, Nan Jiang, Bin Wang, Bo Dai, Puahao Li, Yixin Zhu, Song-Chun Zhu, and Siyuan Huang. 2024. PhyRecon: Physically Plausible Neural Scene Reconstruction. *Advances in Neural Information Processing Systems*.
- Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaeldin El-Nouby, et al. 2023. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193* (2023).
- Luigi Piccinelli, Yung-Hsu Yang, Christos Sakaridis, Mattia Segu, Siyuan Li, Luc Van Gool, and Fisher Yu. 2024. UniDepth: Universal Monocular Metric Depth Estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10106–10116.
- Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. 2022. Dreamfusion: Text-to-3d using 2d diffusion. In *arXiv preprint arXiv:2209.14988*.
- Krishan Rana, Jesse Haviland, Sourav Garg, Jad Abou-Chakra, Ian Reid, and Niko Suenderhauf. 2023. Sayplan: Grounding large language models using 3d scene graphs for scalable robot task planning. In *7th Annual Conference on Robot Learning*.
- Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. 2024. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714* (2024).
- Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kunchang Li, He Cao, Jiayu Chen, Xinyu Huang, Yukang Chen, Feng Yan, et al. 2024. Grounded sam: Assembling open-world models for diverse visual tasks. *arXiv preprint arXiv:2401.14159* (2024).
- Xingyuan Sun, Jiajun Wu, Xiuming Zhang, Zhoutong Zhang, Chengkai Zhang, Tianfan Xue, Joshua B Tenenbaum, and William T Freeman. 2018. Pix3d: Dataset and methods for single-image 3d shape modeling. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2974–2983.
- Stanislaw Szymanowicz, Eldar Insafutdinov, Chuanxia Zheng, Dylan Campbell, João F Henriques, Christian Rupprecht, and Andrea Vedaldi. 2024a. Flash3D: Feed-Forward Generalisable 3D Scene Reconstruction from a Single Image. *arXiv preprint arXiv:2406.04343* (2024).
- Stanislaw Szymanowicz, Chrisitian Rupprecht, and Andrea Vedaldi. 2024b. Splatter image: Ultra-fast single-view 3d reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10208–10217.
- Jiaxiang Tang, Zhaoxi Chen, Xiaokang Chen, Tengfei Wang, Gang Zeng, and Ziwei Liu. 2025. Lgm: Large multi-view gaussian model for high-resolution 3d content creation. In *European Conference on Computer Vision*. Springer, 1–18.
- Jiaxiang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. 2023. Dreamgaussian: Generative gaussian splatting for efficient 3d content creation. *arXiv preprint arXiv:2309.16653* (2023).
- Fengrui Tian, Shaoyi Du, and Yueqi Duan. 2023. Mononerf: Learning a generalizable dynamic radiance field from monocular videos. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 17903–17913.
- Marcel Torne, Anthony Simeonov, Zechu Li, April Chan, Tao Chen, Abhishek Gupta, and Pulkit Agrawal. 2024. Reconciling reality through simulation: A real-to-sim-to-real approach for robust manipulation. *arXiv preprint arXiv:2403.03949* (2024).
- Shinji Ueyama. 1991. Least-squares estimation of transformation parameters between two point patterns. *IEEE Transactions on Pattern Analysis & Machine Intelligence* 13, 04 (1991), 376–380.
- Vikram Voleti, Chun-Han Yao, Mark Boss, Adam Letts, David Pankratz, Dmitry Tochilkin, Christian Laforte, Robin Rombach, and Varun Jampani. 2025. Sv3d: Novel multi-view synthesis and 3d generation from a single image using latent video diffusion. In *European Conference on Computer Vision*. Springer, 439–457.
- Ruicheng Wang, Sicheng Xu, Cassie Dai, Jianfeng Xiang, Yu Deng, Xin Tong, and Jiaolong Yang. 2024b. Moge: Unlocking accurate monocular geometry estimation for open-domain images with optimal training supervision. *arXiv preprint arXiv:2410.19115* (2024).
- Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. 2024a. Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation. *Advances in Neural Information Processing Systems* 36 (2024).
- Xinyue Wei, Minghua Liu, Zhan Ling, and Hao Su. 2022. Approximate convex decomposition for 3d meshes with collision-aware concavity and tree search. *ACM Transactions on Graphics (TOG)* 41, 4 (2022), 1–18.
- Kailu Wu, Fangfu Liu, Zhihan Cai, Runjie Yan, Hanyang Wang, Yating Hu, Yueqi Duan, and Kaisheng Ma. 2024a. Unique3D: High-Quality and Efficient 3D Mesh Generation from a Single Image. *arXiv preprint arXiv:2405.20343* (2024).
- Rundi Wu, Ben Mildenhall, Philipp Henzler, Keunhong Park, Ruiqi Gao, Daniel Watson, Pratul P Srinivasan, Dor Verbin, Jonathan T Barron, Ben Poole, et al. 2024b. Reconfusion: 3d reconstruction with diffusion priors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 21551–21561.
- 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).
- Bin Xiao, Haiping Wu, Weijian Xu, Xiyang Dai, Houdong Hu, Yumao Lu, Michael Zeng, Ce Liu, and Lu Yuan. 2024. Florence-2: Advancing a unified representation for a variety of vision tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4818–4829.
- Tianyi Xie, Zeshun Zong, Yuxing Qiu, Xuan Li, Yutao Feng, Yin Yang, and Chenfanfu Jiang. 2024. Physgaussian: Physics-integrated 3d gaussians for generative dynamics. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4389–4398.
- Qingshan Xu, Jiao Liu, Melvin Wong, Caishun Chen, and Yew-Soon Ong. 2024. Precise-Physics Driven Text-to-3D Generation. *arXiv preprint arXiv:2403.12438* (2024).
- Jianwei Yang, Hao Zhang, Feng Li, Xueyan Zou, Chunyuan Li, and Jianfeng Gao. 2023. Set-of-Mark Prompting Unleashes Extraordinary Visual Grounding in GPT-4V. *arXiv preprint arXiv:2310.11441* (2023).
- Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiaoshi Feng, and Hengshuang Zhao. 2024b. Depth anything: Unleashing the power of large-scale unlabeled data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10371–10381.
- Yandan Yang, Baoxiong Jia, Peiyuan Zhi, and Siyuan Huang. 2024a. Physcene: Physically interactive 3d scene synthesis for embodied ai. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 16262–16272.
- Wei Yin, Chi Zhang, Hao Chen, Zhipeng Cai, Gang Yu, Kaixuan Wang, Xiaozhi Chen, and Chunhua Shen. 2023. Metric3d: Towards zero-shot metric 3d prediction from a

- single image. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 9043–9053.
- Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. 2021. pixelnerf: Neural radiance fields from one or few images. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 4578–4587.
- Hong-Xing Yu, Haoyi Duan, Junhwa Hur, Kyle Sargent, Michael Rubinstein, William T Freeman, Forrester Cole, Deqing Sun, Noah Snavely, Jiajun Wu, et al. 2024. Wonderjourney: Going from anywhere to everywhere. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6658–6667.
- Zehao Yu, Songyou Peng, Michael Niemeyer, Torsten Sattler, and Andreas Geiger. 2022. Monosdf: Exploring monocular geometric cues for neural implicit surface reconstruction. *Advances in neural information processing systems* 35 (2022), 25018–25032.
- Alan Yuille and Daniel Kersten. 2006. Vision as Bayesian inference: analysis by synthesis? *Trends in cognitive sciences* 10, 7 (2006), 301–308.
- Biao Zhang, Jiapeng Tang, Matthias Niessner, and Peter Wonka. 2023. 3dshape2vecset: A 3d shape representation for neural fields and generative diffusion models. *ACM Transactions on Graphics (TOG)* 42, 4 (2023), 1–16.
- Jiancheng Zhang, Haijin Zeng, Yongyong Chen, Dengxiu Yu, and Yin-Ping Zhao. 2024b. Improving Spectral Snapshot Reconstruction with Spectral-Spatial Rectification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 25817–25826.
- Longwen Zhang, Ziyu Wang, Qixuan Zhang, Qiwei Qiu, Anqi Pang, Haoran Jiang, Wei Yang, Lan Xu, and Jingyi Yu. 2024a. CLAY: A Controllable Large-scale Generative Model for Creating High-quality 3D Assets. *ACM Transactions on Graphics (TOG)* 43, 4 (2024), 1–20.
- SUN Zhengwentai. 2023. clip-score: CLIP Score for PyTorch. <https://github.com/taited/clip-score>. Version 0.1.1.
- Licheng Zhong, Hong-Xing Yu, Jiajun Wu, and Yunzhu Li. 2025. Reconstruction and simulation of elastic objects with spring-mass 3d gaussians. In *European Conference on Computer Vision*. Springer, 407–423.
- Qian-Yi Zhou, Jaesik Park, and Vladlen Koltun. 2018. Open3D: A Modern Library for 3D Data Processing. *ArXiv* abs/1801.09847 (2018).

## A GPT-4V PROMPT

The following prompt instructs GPT-4V to act as an object–relationship analyst for a numbered scene image. It defines six relationship types (Stack, Lean, Hang, Clamped, Contained, Edge/Point), sets strict formatting rules (JSON objects in a list). For full details, see the appendix prompt as below.

Listing 1. System prompt for GPT-4V

```

prompting_text_system = """
You are an expert in object recognition and
spatial reasoning.

### Task Description ###
Analyze an image with numbered objects and
determine their relationships.
For each pair of related objects, output a JSON
object containing the relationship details.
Ensure you output all possible relationships, even
those that may be difficult to judge or less
obvious.

### Relationship Definitions ###
1. **Stack**: Object 1 is on top of Object 2 (
   Object 2 supports Object 1 from below)
2. **Lean**: Object 1 is leaning against Object 2
   (Object 2 supports Object 1 laterally)
3. **Hang**: Object 1 is hanging from Object 2 (
   Object 2 supports Object 1 from above)
4. **Clamped**: Object 1 is clamped by Object 2 (
   Object 2 grips Object 1 on multiple sides)
5. **Contained**: Object 1 is inside Object 2 (
   Object 2 encloses Object 1)
6. **Edge/Point**: Object 1 is touching Object 2
   at an edge or point (minimal contact, no
   significant support)

### Important Note ###
Only objects that are in contact with each other
should have a relationship.
For each relationship, always ensure the following
:
1. Use the correct relationship type.
2. Provide a clear explanation of the relationship
.
3. For cases that are in contact but hard to
choose which type, use "Stack".
### Output Format ###
{
  'pair': [obj1_num, obj2_num],
  'relationship': 'Stack'/'Lean'/'Hang'/'Clamped'
    '/Contained'/'Edge/Point',
  'reason': 'explanation'
}

### Examples ###
{ 'pair': [1, 2], 'relationship': 'Stack',      '
  reason': 'Book (1) is stacked on top of Table
(2)' }
```

```

{ 'pair': [3, 4], 'relationship': 'Lean',      '
  reason': 'A chair (3) is leaning against a
  wall (4)' }
{ 'pair': [5, 6], 'relationship': 'Hang',      '
  reason': 'A lamp (5) is hanging from the
  ceiling (6)' }
{ 'pair': [7, 8], 'relationship': 'Clamped',   '
  reason': 'A pipe (7) is clamped by a bracket
  (8)' }
{ 'pair': [9, 10], 'relationship': 'Contained', '
  reason': 'A pencil (9) is inside a pencil case
  (10)' }
{ 'pair': [11, 12], 'relationship': 'Edge/Point', '
  reason': 'A book (11) and a pen (12) are
  touching at the edge' }
Only objects that are in contact with each other
should have a relationship.
"""

content_user = [
  {
    "type": "text",
    "text": """
      ### Object Details ###
      "I have labeled a bright numeric ID at
      the center for each visual object
      in the image."
      "Please analyze all relationships
      between the numbered objects and
      output JSON objects following the
      specified format. "
      "Ensure each relationship includes:"
      "1. The correct relationship type"
      "2. A clear reason for the
      relationship"
      "List all relationships as a JSON
      array."
    """),
  },
  {
    "type": "image_url",
    "image_url": {"url": f"data:image/png;
      base64,{base64_image}"}
  },
]
```