

# Scenethesis: A Language and Vision Agentic Framework for 3D Scene Generation

Lu Ling<sup>1,2</sup>, Chen-Hsuan Lin<sup>1</sup>, Tsung-Yi Lin<sup>1</sup>, Yifan Ding<sup>1</sup>, Yu Zeng<sup>1</sup>, Yichen Sheng<sup>1</sup>, Yunhao Ge<sup>1</sup>, Ming-Yu Liu<sup>1</sup>, Aniket Bera<sup>2\*</sup>, Zhaoshuo Li<sup>1\*</sup>

<sup>1</sup>NVIDIA Research    <sup>2</sup>Purdue University

<https://research.nvidia.com/labs/dir/scenethesis>



Figure 1. *Scenethesis* is a framework for text to interactive 3D scene generation. Given a text prompt, *Scenethesis* leverages both language and visual priors to generate realistic and physical plausible indoor and outdoor environments.

## Abstract

Synthesizing interactive 3D scenes from text is essential for gaming, virtual reality, and embodied AI. However, existing methods face several challenges. Learning-based approaches depend on small-scale indoor datasets, limiting the scene diversity and layout complexity. While large language models (LLMs) can leverage diverse text-domain knowledge,

they struggle with spatial realism, often producing unnatural object placements that fail to respect common sense. Our key insight is that **vision perception can bridge this gap** by providing realistic spatial guidance that LLMs lack. To this end, we introduce *Scenethesis*, a training-free agentic framework that integrates **LLM-based scene planning with vision-guided layout refinement**. Given a text prompt, *Scenethesis* first employs an LLM to draft a coarse layout. A vision module then refines it by generating an image guidance and extracting scene structure to capture inter-object

\*Co-last author.

relations. Next, an optimization module iteratively enforces accurate pose alignment and physical plausibility, preventing artifacts like object penetration and instability. Finally, a judge module verifies spatial coherence. Comprehensive experiments show that *Scenethesis* generates diverse, realistic, and physically plausible 3D interactive scenes, making it valuable for virtual content creation, simulation environments, and embodied AI research.

## 1. Introduction

Synthesizing interactive 3D scenes from text is crucial for gaming [15], virtual content creation [33], and embodied AI [7, 18, 19, 32, 47, 48]. Instead of generating a single scene geometry [14] or differentiable rendering primitives [49], interactive 3D scene synthesis focuses on arranging individual objects to construct a realistic layout while preserving natural interactions, function roles, and physical principles. For example, chairs should face tables to accommodate seating, and small items are typically placed inside cabinets, drawers, and shelves without penetration. Capturing these spatial relationships is crucial for generating realistic scenes, allowing virtual environments to reflect real-world structure and coherence.

Traditional interactive scene generation methods, including manual design [13, 18, 21], are often labor intensive and thus unscalable, while procedural approaches [6] produce overly simplified scenes and fail to capture various real-world spatial relations. In recent years, deep learning-based scene generation methods, such as auto-regressive models [34] and diffusion approaches [41, 47], have enabled end-to-end generation of 3D layouts. However, they rely on object-annotated datasets like 3D-FRONT [12], which are small in scale, limited to indoor environments, and often contain collisions [47]. These datasets primarily model large furniture layouts while neglecting smaller objects and their functional interactions.

The emergence of large language models (LLMs) [10, 20, 48] expands scene diversity by leveraging common-sense knowledge from text, such as which objects should co-occur based on human intent. However, their lack of visual perception prevents them from accurately reproducing real-world spatial relations, leading to unrealistic object placements that disregard functional roles, human intent, and physical constraints. As illustrated in Figure 2, LLM-generated scenes often misorient (e.g., chairs facing the cabinet) and misplace (e.g., cabinet placed against the window) objects; small objects are restricted to predefined locations (e.g., only on top of cabinets instead of inside). This lack of realism disrupts object functionality, weakens spatial coherence, and hinders structural consistency, ultimately making LLM-generated scenes impractical for real-world usability and interactions.

Building on insights from vision foundation models that



Figure 2. Unrealistic 3D scenes generated by the LLM-based method (Holodeck [48]), exhibiting misplaced objects and oversimplified spatial relations.

encode compact spatial information and generate coherent scene distributions reflecting real-world layouts, we introduce *Scenethesis* – a training-free agentic framework that integrates LLM-based scene planning with vision-guided spatial refinement. Building on top of LLMs, which lack real-world perception, *Scenethesis* enforces vision-based spatial constraints to enhance realism and physical plausibility. Given a text prompt, *Scenethesis* employs an LLM for reasoning of coarse layout, a vision module for layout refinement, depth estimation, structural extraction, and a novel optimization for iterative alignment of object placement with visual prior through semantic correspondence matching and signed distance field (SDF)-based physical constraints, ensuring collision-free and stable integration into digital environments. Finally, a judge module verifies the spatial coherence. Quantitative and qualitative results demonstrate that *Scenethesis* outperforms SOTA methods in scene diversity (generating indoor and outdoor scenes), layout realism, and physical plausibility. The layouts generated from *Scenethesis* can be used for downstream tasks such as virtual content creation, editing, and simulation. Our contribution is summarized as follows.

- We introduce *Scenethesis*, a training-free agentic framework, integrates LLMs, vision foundation models, physical-aware optimization, and scene judgment to collaboratively generate realistic 3D interactive scenes.
- *Scenethesis* integrates LLM’s common-sense reasoning for coarse scene planning with vision-guided spatial refinement, effectively capturing realistic inter-object relations.
- We propose a novel optimization process that iteratively aligns objects using semantic correspondence matching and SDF-based physical constraints, enforcing collision-free, stable, and semantically correct placements.
- We assess the diversity, layout realism, and object interactability of scenes generated by *Scenethesis*, demonstrating superior spatial realism and physical plausibility compared to SOTA methods.

## 2. Related Work

**Indoor Scene Synthesis.** Realistic indoor scene synthesis is essential for simulating interactive environments

and training embodied agents for real-world tasks. Early methods framed this task as layout prediction, representing scenes as graphs with object relations [3, 30, 56] or hierarchical structures [23, 42]. SceneFormer [43] and ATISS [34] introduced autoregressive models to infer spatial relations with 3D bounding box supervision. Recent approaches learn layout distributions from 3D datasets like 3D-FRONT [12], while DiffuScene [41] and InstructScene [25] integrate object semantics and geometry into diffusion processes. PhyScene [47] incorporates physical constraints. However, interactive scene generation methods remain dataset-constrained, limiting generalization and often producing unrealistic compositions due to relaxed collision constraints [41, 43]. Instead of learning layout distributions from limited 3D datasets, *Scenethesis* derives spatial priors from image generation models, enabling broader generalization across diverse scenarios.

**LLM/VLM Guided 3D Scene Generation.** Early efforts [6, 7, 37] relied on rule-based procedural modeling to define spatial relations for interactive environments. With the rise of LLMs/VLMs, recent methods such as SceneTeller [33], Holodeck [48], SceneCraft [20], GALA3D [55], RobotGen [44], Open-Universe [1], GenUSD [26], LayoutVLM [40] and SceneX [54] leverage LLMs/VLMs for: (1) spatial relation planning via predefined implicit relations, (2) 3D asset retrieval from semantic descriptions or vision-language embeddings, and (3) rule-based rough collision detection, demonstrating large-scale scene generation potential. Although LLMs encode rich common sense knowledge, they struggle with fine-grained spatial reasoning. Predefined spatial relations in text descriptions are often simplistic, limiting their ability to capture the complexity of the real-world scene [17, 25]. In contrast, *Scenethesis* leverages LLM priors to convert text prompts into coarse layout instructions while using vision foundation model to persevere compact spatial information, effectively capturing real-world spatial complexity.

**Visual Foundation Model-Guided Scene Generation.** Visual foundation models (VFs), particularly image generation models, have advanced visual generation and are now widely applied to 3D scene synthesis. Methods such as Text2Room [14], SceneScape [11], WonderJourney [50], WonderWorld [49], and Text2NeRF [52] integrate 2D diffusion with 3D priors (e.g., depth) to generate single-geometry scenes. However, this approach inherently faces challenges in handling occlusions and reconstructing hidden elements due to the interconnected structure of real-world scenes, making them unsuitable for object interactions.

Architect [45] and Deep Prior Assembly (DPA) [53] introduce 2D inpainting for interactive 3D scene generation and reconstruction. While this improves occlusion handling, the lack of physical constraints and 3D reasoning leads to misaligned, floating, or intersecting objects, making it diffi-

cult to maintain functional object relationships for embodied AI tasks. In contrast, *Scenethesis* integrates physics-aware optimization, ensuring both spatial alignment with realistic visual prior and physical plausibility.

**Physics-Aware Scene Generation.** Physical principles have been largely overlooked in 3D interactive scene generation for both LLM-based and VFM-based methods. Recent works, such as PhyScene [47] and Holodeck [48] enforce physical constraints by detecting collisions using 3D bounding boxes. While PhyScene reduces collision rates, it still exceeds 15% [47]. Holodeck focuses only on large-object collision avoidance, neglecting small-object inter-collisions. Despite these advances, achieving full physical plausibility remains a challenge. To address this, *Scenethesis* incorporates precise collision detection and stability constraints, significantly reducing collision and instability rates.

### 3. Method

*Scenethesis* generates **spatially realistic, physically plausible** interactive 3D environments from user prompts. An overview of the pipeline is shown in Figure 3, consisting of four key stages: (1) an **LLM module** drafts a coarse scene plan, (2) a **vision module** refines the layout with visual guidance and structural extraction, (3) a **physical-aware optimization module** distills priors and adjusts object placement for spatial coherence and physical plausibility, and (4) a **scene judge module** verifies spatial consistency. The following sections detail each module’s role.

#### 3.1. Coarse Scene Planning

*Scenethesis* supports either a **simple prompt** (e.g., “a peaceful beach during sunset”) for flexible scene generation or a **detailed prompt** for controllable scene generation (e.g., a scene plan describing the detailed spatial relations as shown in the *appendix*). For a simple prompt, the **LLM** generates a coarse scene plan by reasoning over user input. It first interprets the prompt, reviews all object categories in the available 3D database, selects commonly associated objects, and then generates an up-sampled prompt describing coarse spatial relations, as illustrated in Figure 3. When given detailed prompts, the LLM checks for the presence of all specified objects in the database, infers relevant object categories, and skips the prompt up-sampling process.

Among the selected objects, the **LLM** identifies an *anchor* object, following prior work [48]. The anchor serves as the central reference point, occupying the highest spatial hierarchy apart from the *ground*. Then the LLM establishes a coarse spatial hierarchy, positioning objects relative to the anchor and incorporating these relationships into the upsampled prompt. For example, in a cozy living room, the sofa acts as the anchor at the *center*, while a bookshelf is placed in the *background, aligned against* the wall. Other objects,

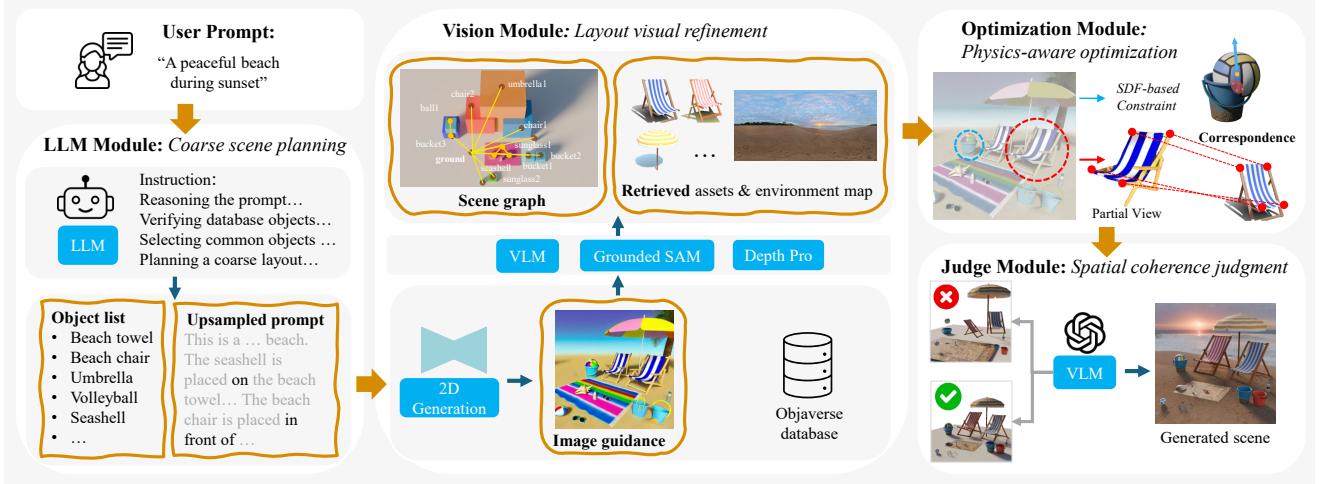


Figure 3. *Scenethesis* is an agentic framework. The LLM module performs coarse scene planning, estimating rough spatial relationships. The vision module refines this layout by enforcing accurate spatial constraints. The physical-aware optimization iteratively adjusts object placement, ensuring pose alignment and physical plausibility. Finally, a judge module verifies the scene spatial coherence.

such as a coffee table or chairs, are positioned *in front of* or *beside* the sofa.

### 3.2. Layout Visual Refinement

A key insight of *Scenethesis* is that image generation models inherently encode object functionality and spatial relationships by learning common co-occurrences and spatial arrangements from large-scale image datasets. The **vision module** refines the coarse layout through three steps: (1) *Image Guidance* – Generates images to refine spatial relations, ensuring realism and object functionality. (2) *Scene Graph Generation* – Segments objects, estimates depth and 3D bounding boxes, and constructs a graph encoding inter-object relationships to establish the initial layout. (3) *Asset Retrieval* – Selects 3D assets and environment maps for final scene composition.

**Image Generation.** The **vision module** refines the upsampled prompt into a visually structured scene representation. This generated image serves as the basis for segmentation, depth estimation, and asset retrieval.

**Scene Graph Generation.** Leveraging vision foundation models such as GPT-4o [16], Grounded-SAM [39], and DepthPro [2], the vision module constructs a **scene graph** that localizes objects using 3D bounding boxes (3DBB) and identifies structural components, including the *anchor object*, *parent objects*, and *child objects* (see Figure 3).

To initialize asset 5DoF poses, vision module segments objects using semantic cues, estimates depth maps, and projects them into a 3D point cloud. However, due to occlusion, limited perspectives, and segmentation errors, cropped image guidance may miss full object visibility, leading to biases in 3DBB estimation – necessitating pose adjustments later (Sec. 3.3.1).

The scene graph forms the basis for iterative 5DoF pose

adjustments during optimization in the next stage. Since *Scenethesis* focuses on ground-level object layout, background elements, e.g. wall decorations, are visually defined by the retrieved environment map. Detailed scene graph formatting instructions are provided in the *appendix*.

**Asset Retrieval.** Unlike existing 3D object generation and reconstruction techniques [24, 46], such as 3D Gaussian splatting, which can produce photorealistic visuals but suffer from artifacts and geometric inconsistencies. These methods lack editable meshes, UV mappings, and decomposable PBR materials, making them incompatible with standard production workflows. To address these limitations, *Scenethesis* adopts a retrieval-based approach for asset selection, ensuring both geometric fidelity and editability for downstream applications. We construct a high-quality asset subset from Objaverse [8] similar to Holodeck [48], and supplemented with a custom environment map dataset. In the final step, the 3D assets and an environment map are retrieved to assemble a visually coherent scene. Retrieval details can be found in *appendix*.

### 3.3. Physics-aware Optimization

Directly placing 3D assets based on estimated point clouds from image guidance poses significant challenges: (1) *Occlusions* in real-world scenarios result in incomplete 3D point clouds, leading to errors in object orientation, scale, and position. (2) *Discrepancies* between retrieved assets and image guidance in texture and shape make precise pose estimation difficult. To overcome these issues, *Scenethesis* employs a **physics-aware optimization** powered by robust semantic feature matching [4, 9, 51] and signed-distance fields (SDFs). This optimization process iteratively refines object poses to ensure pose alignment and physical plausibility.

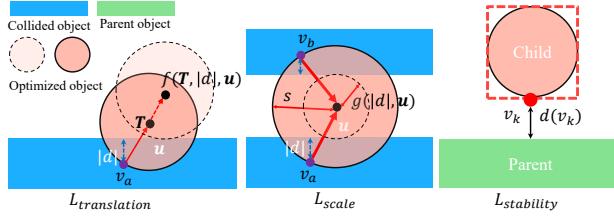


Figure 4. Collision avoidance and stability maintenance.

### 3.3.1. Pose Alignment

To address pose estimation errors from occlusions, segmentation, or asset mismatches, we adopt dense correspondence matching from RoMa [9], leveraging semantic spatial features for robustness to occlusions and partial views. Unavoidable discrepancies in texture and shape between image guidance and retrieved assets are mitigated by focusing on high-level semantics over low-level details.

For each object, we match  $N$  correspondences between the rendered object and partially visible regions in the image guidance in 2D space. It then minimizes MSE loss on both 2D and 3D spatial locations of these  $N$  correspondences, backpropagating gradients to refine scale, translation, and upright rotation, as shown in Figure 3. Further details on pose estimation are provided in the Appendix.

### 3.3.2. Physical Plausibility

Real-world 3D scenes obey physical constraints, ensuring objects remain stable on contact surfaces and collision-free. However, pose alignment with image guidance alone does not guarantee physical plausibility—objects may intersect, float, or sink due to shape discrepancies and errors in scene understanding. See Figure 9 (b) as an example.

Existing methods approximate object geometry using 3D bounding boxes (3DBB) [47, 48], which oversimplifies shapes and leads to simplified inter-object relationships. For example, objects cannot be put within the shelf due to 3D bounding box collision. This results in simplified scene diversity, especially in tight spaces with complex inter-object relationships (see Figure 8 for an example). To address these challenges, we replace 3DBB-based approximations with Signed Distance Fields (SDFs), enabling precise object geometry representation for accurate collision detection and stability constraints.

The **physical-aware optimization process** iteratively constructs a SDF-based physical structure, following the scene graph hierarchy: processing the anchor object first to establish a stable foundation, followed by parent and child objects. The physics-aware optimization incorporates **collision** and **stability constraints**. Since retrieved 3D assets are upright, their rotation is constrained to azimuthal adjustments.

Formally, given a scene graph with  $N$  objects, each object has a 5-DoF configuration defined by scale  $s$ , upright rotation  $\mathbf{R}$ , and translation  $\mathbf{T} = (t_x, t_y, t_z)$ . For computational

efficiency, we uniformly sample  $n$  points from its triangle surface mesh as its geometric representation and compute its centroid for collision avoidance.

**Collision Constraints.** We query the scene SDFs using object surface points to detect collision states and define position collision loss  $\mathcal{L}_{\text{translation}}$  and scale collision loss  $\mathcal{L}_{\text{scale}}$ . As shown in Figure 4, the deviation caused by collisions impacts translation  $T$  as:

$$\mathcal{L}_{\text{translation}} = \sum_{\mathbf{v}_i \in \mathbf{V}^-} \|f(\mathbf{T}, |d_i|, \mathbf{u}_i) - \mathbf{T}\|_2^2, \quad (1)$$

where  $f(\mathbf{T}, |d_i|, \mathbf{u}_i) = \mathbf{T} + \mathbf{u}_i \cdot |d_i|$  computes a collision-free position by adjusting the translation along direction  $\mathbf{u}_i$  with step size  $|d_i|$ . Here,  $d_i$  is the negative SDF value at a collided point  $\mathbf{v}_i$ , which belong to the points set with negative SDF  $\mathbf{V}^-$  sampled uniformly from the surface. The direction  $\mathbf{u}$  is defined from the collision point toward the model’s centroid, guiding objects away from collisions.

Collisions also affect object scale  $s$  due to opposing forces:

$$\mathcal{L}_{\text{scale}} = \begin{cases} \sum_{\mathbf{v}_i \in \mathbf{V}^-} \left( g(|d_i|, \mathbf{u}_i) - s \right)^2 & \text{if } N_{\text{cluster}} > 1, \\ 0 & \text{otherwise} \end{cases}, \quad (2)$$

where  $g(|d_i|, \mathbf{u}_i) = \frac{\|\mathbf{u}_i\| - |d_i|}{\|\mathbf{u}_i\|}$  defines the target scale to reduce collision regions.  $N_{\text{cluster}}$  denotes the number of distinct clusters formed without SDF sign flipping. As shown in Figure 4, two surface points  $i$  and  $j$  with  $d_i \leq 0$  and  $d_j \leq 0$  belong to different clusters, and thus push the object to be smaller.

**Stability Constraints.** Objects are dragged by gravity and rest on their bottom contacting surface. We ensure stability by enforcing contact between an object’s bottom points and its parent surface, where their SDF values should be zero, as shown in Figure 4. The stability loss is defined as:

$$\mathcal{L}_{\text{stability}} = \sum_{\mathbf{v}_i \in \mathbf{V}^B} \left( 1 - \exp(-d_i^2) \right), \quad (3)$$

where  $\mathbf{V}^B$  are the sampled points on the bottom surface of bounding box, and  $d_i$  are their corresponding SDF values. Further details on collision loss optimization are provided in the Appendix.

### 3.4. Spatial Coherence Judgment

After iteratively optimizing object placement, a scene judge powered by GPT-4o evaluates the spatial alignment between the generated 3D scene and the image guidance produced during the layout refinement stage, ensuring consistency in inter-object relationships.

To assess this alignment, we design three metrics: (1) object category accuracy, comparing the generated scene

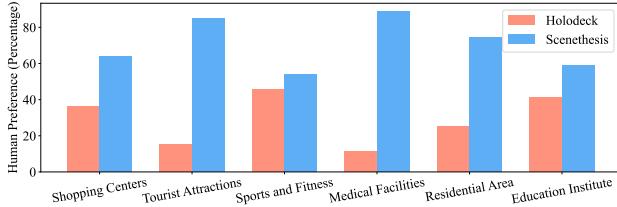


Figure 5. Human preference on diverse indoor scenes.

with the image guidance; (2) object orientation alignment, measuring how well object orientations match the reference layout; (3) overall spatial coherence, capturing the holistic consistency of the scene layout.

Each metric is normalized between 0 (lowest) and 1 (highest). If any metric falls below a predefined threshold, the scene judge triggers a re-planning step. Further details are provided in the *Appendix*.

## 4. Experiment

**Implementation.** We use GPT-4o [16] as the LLM and image generation in vision module. Following Holodeck [48], we retrieve 3D models from a high-quality Objaverse [8] subset. Other module details are discussed in the above section. The physics-aware optimization is implemented by PyTorch [35] and PyTorch3D [38]. Experiments are run on an A100 GPU.

**Baselines.** Since we focus on interactive scene generation, methods producing only single-geometry representations are not relevant. We compare our approach against open-sourced state-of-the-art (SOTA) *generative methods* (DiffuScene [41], PhyScene [47]) and *LLM-based methods* (SceneTeller [33], Holodeck [48]). For fairness, all LLM-based methods use the same ChatGPT version.

**Setup.** *Scenethesis* generates both indoor and outdoor scenes (Figure 1, Figure 6), but for fair comparison, we evaluate only indoor scenes. To assess diversity and realism, we generate 22 indoor scenes covering 6 primary and 12 secondary categories from DL3DV-10K [29]: **Residential** (*living room, playroom, garage, warehouse*), **Shopping** (*bookstore, store*), **Tourism** (*museum, piano showroom*), **Sports** (*gym, billiard club*), **Medical** (*ward*), **Education** (*laboratory*). Since DiffuScene, PhyScene, and SceneTeller were trained on indoor datasets [12] (mainly residential areas), we compare them within this domain. Holodeck, which also retrieves models from Objaverse, supports indoor scene generation, enabling comparisons across all indoor categories. To mitigate view-dependent bias, we render each scene from two perspectives, yielding 44 image pairs. For baselines lacking background generation (e.g. SceneTeller), we render *Scenethesis* outputs without an environment map for a fair comparison.

### 4.1. Metrics

We evaluate *controllability* in text-based scene generation methods and three key properties essential for virtual content generation: *layout realism*, *physical plausibility*, and *interactivity*.

**Controllability.** Ensuring 3D scene generation aligns with input prompts is crucial. We assess this using: (1) *CLIP Score* [36] – cosine similarity between image and text features from CLIP. (2) *BLIP Score* [22] – image-text alignment using the ITM head of BLIPv2. (3) *VQA Score* [27] – image-caption alignment based on VQA models.

**Layout Realism.** Visual quality and spatial realism are important to reflect real-world scene layouts. We evaluate it using following metrics: (1) *Object Diversity* – number of objects and categories in the scene. (2) *Layout Coherence* – adherence of object positions and orientations to common sense. (3) *Spatial Realism* – presence of diverse spatial relations (e.g., on top of, inside, under). (4) *Overall Performance* – alignment of object categories and styles with the scene type. Evaluation details and examples are in the *appendix*.

**Physical Plausibility.** Ensuring object collision-free and stable placement is fundamental for physical simulation environments. We construct the following metrics: (1) *Col-O* – average object collision rate, (2) *Col-S* – average scene collision rate, (3) *Inst-O* – average object instability rate, and (4) *Inst-S* – average scene instability rate.

Collision is tested via mesh-mesh intersections, while instability follows Atlas3D [5], measured by tracking transformations after physics-based simulation [31]. These metrics assess scene viability for virtual content creation.

**Interactivity.** To ensure objects are accessible and manipulable in the scene based on their functional roles, we follow evaluation metrics from PhyScene [47]: (1) *Reach* – average object reachability rate, and (2) *Walk* – ratio of the largest connected walkable area over all walkable regions.

### 4.2. Quantitative Evaluation

**Controllability.** Table 1 presents a comprehensive evaluation of text-image alignment. Among all baselines, *Scenethesis* achieves the highest CLIP, BLIP, and VQA scores, confirming its effectiveness in adhering to text description and the reliability of our agentic pipeline.

**Layout Realism.** Table 1 reports visual realism scores from human evaluations and GPT-4o, a human-aligned evaluator in text-to-3D tasks [16, 28]. *Scenethesis* achieves SOTA performance on most metrics. Despite DiffuScene and PhyScene being trained on dedicated indoor residential datasets [12], the training-free *Scenethesis* achieves comparable or superior layout realism in residential areas. In broader indoor settings (e.g., shopping centers, tourist attractions, sports facilities), Table 1 and Figure 5 show that *Scenethesis* significantly outperforms Holodeck in visual quality and spatial realism. These results demonstrate the advantages of

Table 1. Quantitative evaluation on text–image alignment and visual-quality preference ( $\uparrow$  higher is better). **Bold** marks the best for text control measurement. Visual quality preference indicates GPT-4o and human preference for our method over the baseline.

Method	Text–Image Alignment			Visual-Quality Preference of Ours (GPT-4o / Human Evaluation)			
	CLIP $\uparrow$	BLIP $\uparrow$	VQA $\uparrow$	Object Diversity $\uparrow$	Layout Coherence $\uparrow$	Spatial Realism $\uparrow$	Overall Performance $\uparrow$
PhyScene	–	–	–	80% / 75%	60% / 46%	85% / 74%	50% / 53%
DiffuScene	23.11	48.28	0.7832	75% / 80%	80% / 90%	90% / 76%	80% / 80%
SceneTeller	25.27	51.99	0.7999	80% / 85%	80% / 71%	85% / 80%	80% / 74%
Holodeck	28.32	46.25	0.6815	85% / 80%	83% / 78%	81% / 86%	85% / 85%
<b>Ours</b>	<b>30.71</b>	<b>77.17</b>	<b>0.8269</b>	– / –	– / –	– / –	– / –

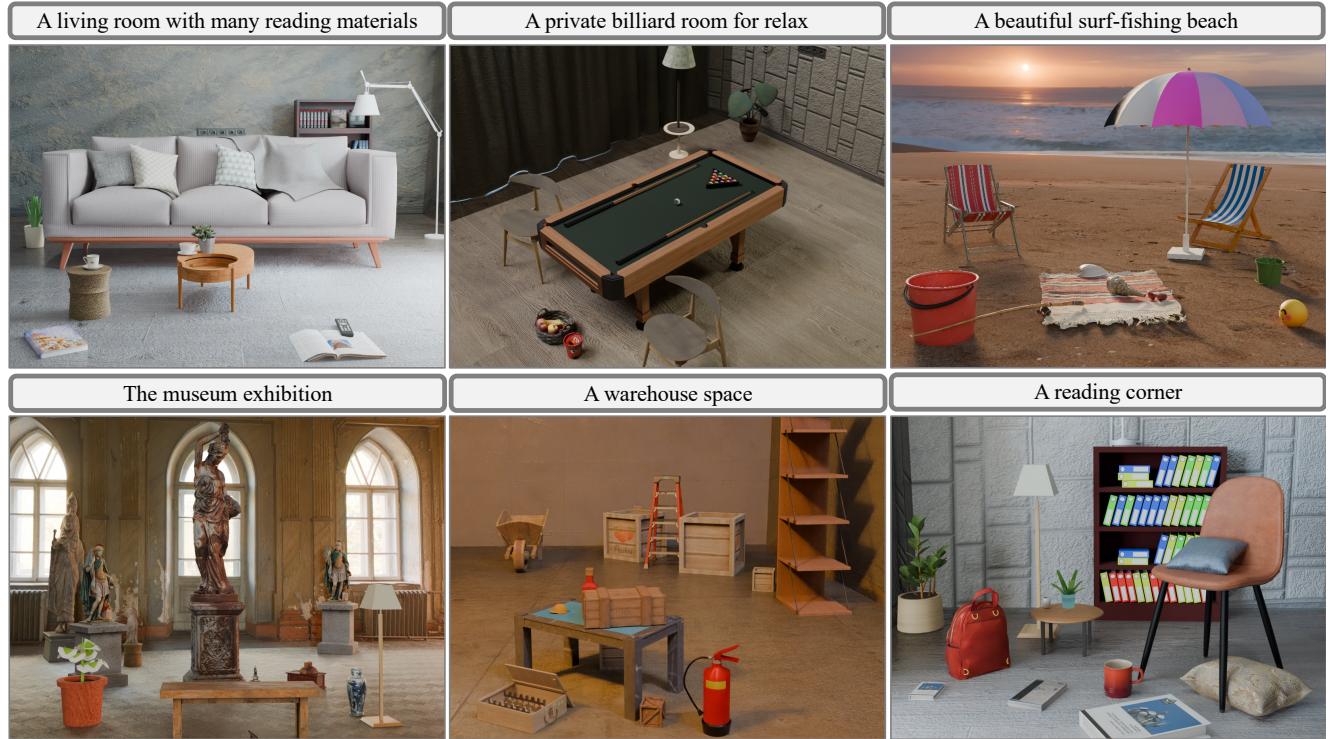


Figure 6. **Qualitative results of generated indoor and outdoor scenes by *Scenethesis*.** *Scenethesis* can generate diverse scenes given user prompts. Visualizations of the scenes at different camera viewpoints can be found in *appendix*.

Table 2. Physical-plausibility and interactivity results ( $\downarrow$  lower is better for collision/instability). **Bold** indicates the best value.

Method	Physical Plausibility				Interactivity	
	Col-O $\downarrow$	Col-S $\downarrow$	Inst-O $\downarrow$	Inst-S $\downarrow$	Reach $\uparrow$	Walk $\uparrow$
PhyScene	17.6%	51%	18.73%	75.22%	0.77	0.84
DiffuScene	19.5%	55%	20.75%	83.33%	0.74	0.83
SceneTeller	35.2%	75%	41.17%	78.57%	0.75	0.80
Holodeck	6.1%	21%	7.00%	31.58%	0.90	<b>0.96</b>
<b>Ours</b>	<b>0.8%</b>	<b>6%</b>	<b>3.20%</b>	<b>16.67%</b>	<b>0.94</b>	<b>0.96</b>

visual prior in guiding spatially realistic scene generation.

**Physical Plausibility and Interactivity.** Table 2 presents object-level and scene-level physical plausibility metrics, demonstrating that *Scenethesis* significantly reduces collisions and enhances stability.

DiffuScene [41] and SceneTeller [33], trained on high-collision datasets [12, 47], lack collision detection and

stability constraints, leading to frequent object intersections. PhyScene [47] applies physical constraints but inherits dataset-induced collisions. Holodeck [48] prevents large-object collisions via Depth-First-Search solver but places small objects on predefined surfaces without collision checks, often causing inter-object penetrations (see *appendix*). Moreover, none of these baselines address stability, resulting in frequent failures in physics-based simulations.

In contrast, *Scenethesis* integrates physics-aware layout adjustment, ensuring low-collision, stable environments. Beyond physical plausibility, *Scenethesis* excels in interactivity, achieving superior reachability and walkability scores. These results highlight *Scenethesis*’s ability to generate accessible, navigable environments where objects align with their functional roles and afford interactions.

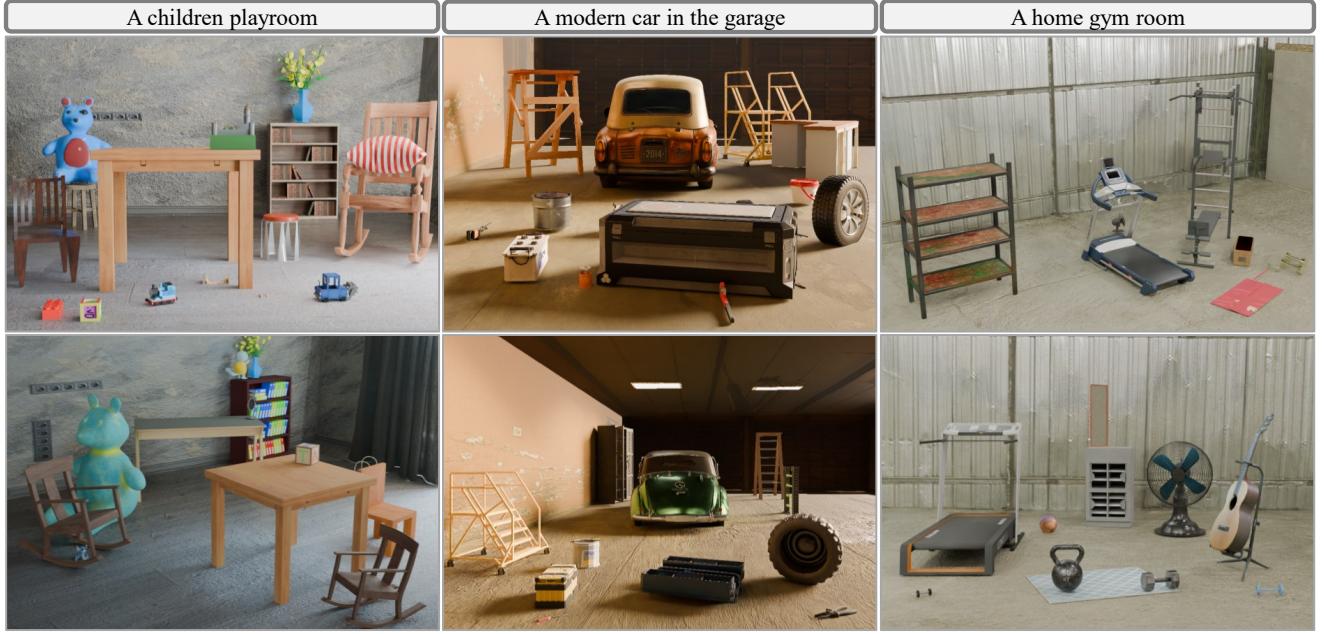


Figure 7. **Output Diversity.** Given the same text prompt, *Scenethesis* can generate diverse scene with various objects and different layouts.



Figure 8. **Complex spatial realism.** Spatial realism comparison between *Scenethesis* and *Holodeck*. *Scenethesis* generates spatially plausible 3D scenes, precisely placing small objects (e.g., bag, wine bottle, vase) within shelf compartments rather than just on top. This precision, challenging for LLM-based methods, is essential for embodied agent manipulation tasks [32, 47].

### 4.3. Qualitative Evaluation

Figure 6 showcases diverse scenes generated by *Scenethesis*, demonstrating high fidelity and versatility in both indoor and outdoor environments. Compared to LLM-based approaches, *Scenethesis* excels in realism and physical plausibility by leveraging image guidance and physics-aware optimization, effectively capturing real-world spatial complexity and diversity. Figure 7 presents various 3D layouts generated from the same text prompt, highlighting diverse asset selection and spatial arrangements. *Scenethesis* supports both simple and detailed prompts—simple prompts enable flexible, user-friendly generation, while detailed prompts allow controllable 3D scene generation (see appendix).

*Holodeck* restricts small object placement to predefined areas on the top of larger objects. In contrast, *Scenethesis* enables fine-grained positioning, placing small object at different levels of the support structure (e.g., shelves, carts), as

Table 3. Ablation study on the effectiveness of physically plausible optimization. *Scenethesis* is the result in “+Stability” which includes all constraint components.

Component	Pose Alignment $\uparrow$	Collision Rate $\downarrow$	Instability Rate $\downarrow$
Raw layout	0.536	22.7%	87.3%
+Pose Alignment	0.732	10.6%	74.2%
+Collision	0.755	3.6%	69.8%
+Stability	<b>0.836</b>	<b>0.8%</b>	<b>3.2%</b>

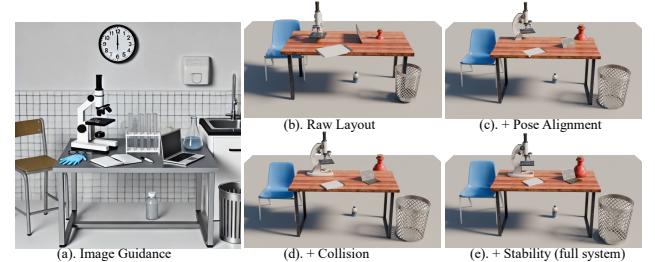


Figure 9. **Effects of different constraints.** (a) *Scenethesis* plans the layout and generates image guidance from text input. (b) *Raw layout*: places 3D models in estimated 3DBBs. (c) + *Pose alignment*: adjusts 5DoF poses but lacks physical plausibility. (d) + *Collision*: prevents intersections but allows floating objects. (e) + *Stability*: ensures grounded, physically stable objects.

shown in Figure 8. LLM-based methods, which lack visual perception, struggle with this level of spatial realism. This capability is critical for embodied AI, enabling realistic interactions and meaningful object manipulation in simulation. More examples and qualitative comparisons can be found in appendix.

#### 4.4. Ablation Study

The physics-aware optimization has three components: *pose alignment*, *collision constraint*, and *stability constraint*. We perform ablation studies to assess their effectiveness.

**Metric.** For each generated scene, we render the same view as the image guidance and use GPT-4o to assess *pose alignment* based on: (1) object orientation, size, and position similarity, and (2) spatial coherence of the overall layout. The similarity score ranges from 0 to 1, with 1 indicating the highest alignment. Wall decorations are ignored in the comparison. Additionally, we evaluate object collisions and instability using the method in [Section 4.2](#).

**Baselines.** *Raw Layout:* Objects are placed based on 3DBB estimated by segmentation and depth prediction methods.

*Pose Alignment:* Aligns object placement with image guidance via correspondence matching. *Collision Constraint:* Optimizes placement to avoid collisions. *Stability Constraint:* Ensures objects remain stable.

**Results.** As shown in [Table 3](#), pose alignment significantly improves spatial consistency, while collision and stability constraints enhance physical plausibility, making scenes simulation-ready. [Figure 9](#) shows qualitative visualization.

### 5. Conclusion and Limitation

We introduce *Scenethesis*, a training-free agentic framework for generating high-fidelity interactive 3D scenes by leveraging LLM-based coarse scene planning, vision-guided layout refinement, and physics-aware optimization for object position adjustment. A scene judge module ensures spatial coherence. Experimental results demonstrate that it significantly outperforms SOTA baselines in layout coherence, spatial realism, and plausibility. Our approach is limited by retrieval databases since generative 3D methods cannot yet handle articulation. Future advances in generative 3D could overcome this constraint by enabling articulated object synthesis, enhancing scene diversity.

# Scenethesis: A Language and Vision Agentic Framework for 3D Scene Generation

## Supplementary Material

### 6. Implementation Details of Scenethesis

#### 6.1. Algorithm Overview

In this section, we provide a high-level algorithmic overview of *Scenethesis*, with detailed steps outlined in Algorithm 1.

#### 6.2. Method Details

##### 6.2.1. Coarse Scene Planning

Using the user’s scene prompt as input, the LLM (powered by GPT-4o [16]) follows a six-step process:

1. Interpreting the user’s scene prompt.
2. Reviewing the object categories available in the provided asset database.
3. Selecting relevant objects from the asset list.
4. Cross-checking the availability of the selected objects.
5. Planning the scene using the selected objects.
6. Generating output files according to the specified standards.

The final coarse scene planning output consists of two components: a list of selected object categories commonly found in the scene (defining anchor object and other common objects) and an upsampled prompt that outlines the scene’s spatial hierarchy. The designed prompt presents in *Coarse Scene Planning Instruction Prompts* Section 7.1 and the output example is in *Coarse Scene Planning Output Example* Section 7.2.

##### 6.2.2. Layout Visual Refinement

Based on the upsampled prompt, GPT-4o generates an image to serve as fine-grained layout guidance. Several post-processing steps are applied to the generated image:

- **Scene Graph Construction:** GPT-4o [16] is used to generate a scene graph, defining the ground as the *root object*, along with *parent objects* and their corresponding *child objects*. Additionally, Grounded-SAM [39] segments each object in the image to obtain masks and cropped images. These are then projected into 3D space using Depth Pro [2], allowing for the initial positioning of objects within a spatial relationship graph.
- **Asset Retrieval.** CLIP (ViT-L/14 trianed on LAION-2B) image and semantic features are employed to retrieve 3D assets that align with the image guidance. GPT-4o [16] is further utilized to select the most relevant environment map based on the upsampled prompt. It is important to note that *Scenethesis* focuses on layout planning for objects on the ground, while background elements, such as wall decorations, lighting, or outdoor settings (e.g., sunshine or the sea), are visually determined by the environment map.

The output of the fine-grained layout planning includes the generated image as guidance, a scene graph with the initial poses of the objects, the retrieved assets, and the retrieved environment map. The visual details are presented in the video.

##### 6.2.3. Physics-aware optimization Details

The physics-aware optimization is an iterative optimization process that consists of two key components: **pose alignment** optimization and **physical plausibility** optimization. pose alignment optimization focuses on aligning the position, dimension, and orientation of 3D models with their counterparts in the image guidance to ensure visual coherence for spatial relationships. Physical plausibility optimization ensures that the 3D models in the scene are free from collisions and maintain stability, contributing to a realistic and physically consistent layout.

**Pose Alignment.** To align the position, dimension, and orientation for objects in rendered image and their counterpart in image guidance, *Scenethesis* applies the dense semantic correspondence matching from RoMa [9]. That is, minimizing the distance between correspondence points in the rendered image  $I$  and the guided image  $\tilde{I}$ . Suppose there are  $N$  objects in the rendered image  $I$ , each represented by  $\mathbf{o}$  and defined by a 5-DoF configuration, which includes scale  $s$ , upright rotation  $\mathbf{R}$ , and translation  $\mathbf{T} = (t_x, t_y, t_z)$ . The counterpart of each object in the generated image  $\tilde{I}$  is denoted as  $\tilde{\mathbf{o}}$ . The objective of ensuring visual coherence is to minimize the distance between corresponding points by optimizing the 5-DoF parameters. This ensures that the spatial positions, dimension, and orientations of the 3D models are closely aligned with their counterparts in the guided image. The matching process is formalized as:

$$\{p(x, y), \tilde{p}(x, y)\}_i^m = \text{RoMa}(\mathbf{o}, \tilde{\mathbf{o}}), \quad (4)$$

where  $p(x, y), \tilde{p}(x, y)$  are correspondent pair in object  $\mathbf{o}$  and  $\tilde{\mathbf{o}}$ . We select  $m$  pair points in each optimization iteration with confident score higher than  $\tau$ . The higher confidence score indicates a higher probability of matching. We minimize the 2D pixel distance and 3D projected point clouds distance between the matched pair denoted as follows:

$$\mathcal{L}_{\text{pose}} = \lambda_{2d}\mathcal{L}_{2d} + \lambda_{3d}\mathcal{L}_{3d}, \quad (5)$$

where  $\lambda_{2d}$  and  $\lambda_{3d}$  are coefficients of the 2D pixel loss and 3D point cloud loss denoted as  $\mathcal{L}_{2d}$  and  $\mathcal{L}_{3d}$ .

**Algorithm 1** Text to 3D Interactive Scene Generation

```

1: Input: User text
2: Output: 3D interactive scene layout
3:
4: Stage Stage 1: Coarse Scene Planning :
5: object_list, upsampled_prompt ← LLM(user_text)           ▷ obtain the object list and an upsampled prompt
6:
7: Stage Stage 2: Layout Visual Refinement :
8: img_guidance ← 2D_Diffusion (upsampled_prompt)
9: cropped_images ← Grounded_SAM (img_guidance, object_list)
10: depth_map ← Depth_Pro (img_guidance)
11: 5DoF_poses ← Extract_Poses(cropped_images, depth_map)
12: scene_graph ← VLM (img_guidance, object_list, 5DoF_poses)
13: 3D_assets← CLIP (cropped_images, object_list)
14: environment_map ← VLM (upsampled_prompt)               ▷ generate the guidance image as the reference
                                                               ▷ identify each object and crop the images
                                                               ▷ generate depth map
                                                               ▷ generate initial 5DoF poses
                                                               ▷ generate scene graph
                                                               ▷ retrieve 3D assets
                                                               ▷ retrieve environment maps
15:
16: Stage Stage 3: Physics-aware Optimization:
17: scene_SDF ← Init_Scene_SDF(anchor_object)             ▷ compute SDF for each object
18: for node in scene_graph.bfs_traverse() do            ▷ iterate over all objects
19:   s, R, T ← node.pose                                ▷ variables to be optimized
20:   parent_SDF ← node.parent.SDF                         ▷ obtain parent object's SDF
21:   for iteration = 1 to max_iterations do
22:     mesh ← Get_Object_Mesh(node)                      ▷ coordinate alignment
23:     mesh* ← Apply_Transform(mesh, s, R, T)
24:     img_rendered, depth_rendered ← Render(mesh*, camera)    ▷ render RGB and depth images
25:     correspondence ← RoMa(img_guidance, img_rendered)      ▷ correspondence matching
26:     mesh_points ← Get_Point_Clouds(depth_rendered, correspondence, camera)
27:     guided_points ← Get_Point_Clouds(depth_map, correspondence, camera)    ▷ loss computation
28:     L_pose_2D ← Dist_2D(correspondence)
29:     L_pose_3D ← Dist_3D(mesh_points, guided_points)
30:     L_collision ← Collision(mesh*, scene_SDF)
31:     L_stability ← Stability(bottom_points(mesh), parent_SDF)
32:     loss ←  $\lambda \mathcal{L}_{pose} + \lambda_{collision} \mathcal{L}_{collision} + \lambda_{stability} \mathcal{L}_{stability}$ 
33:     loss.Backward()                                     ▷ pose optimization
34:     optimizer.Step()
35:     optimizer.Zero_Grad()
36:   end for
37:   scene_SDF ← Update_Scene_SDF(scene_SDF, node)
38: end for
39:
40: Stage Stage 4: Scene Spatial Coherent Judgment:
41: Multi-view images ← Render (optimized_3D_scene)
42: Qualified ← VLM (Multi-view images)
43: if not qualified then
44:   goto Stage 1                                         ▷ re-generate if current optimization fails.
45: end if
46: Return: Optimized 3D interactive scene

```

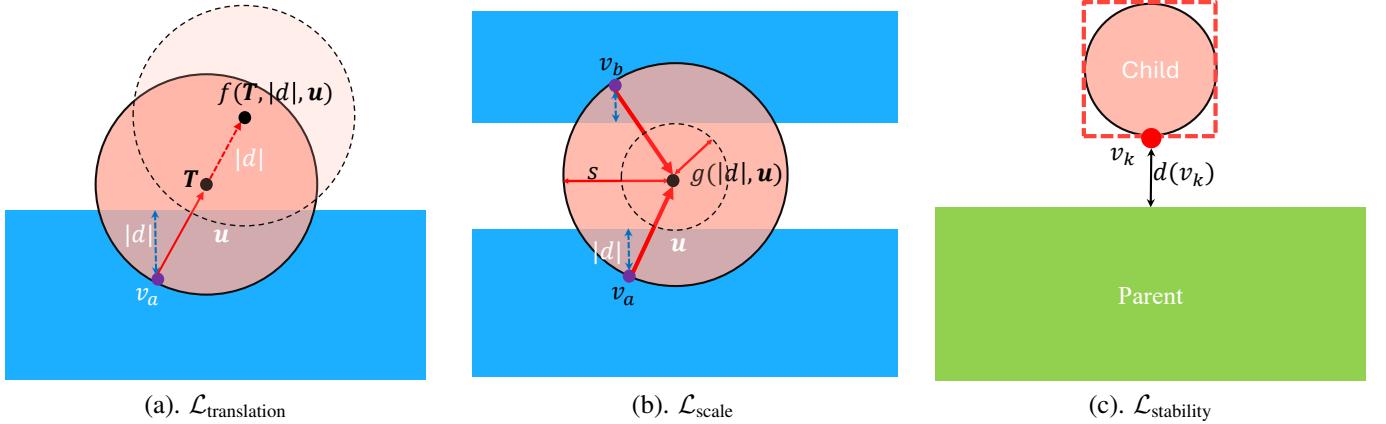


Figure 10. Illustration of *collision avoidance* and *stability maintenance*. The solid-line circle indicates the 3D object’s current position, while the dotted-line circle marks its anticipated position. The black dot represents the centroid of the target object, the purple dots indicate surface nodes with negative SDF values, and the red point  $v_k$  is the bottom node of the object. (a). The collision pushes the circle object out of the rectangle along the direction from the sampled point to the circle’s center by step  $|d|$ . (b). The collision indicates the object is too large and negative signed distance fields (SDF) points (*i.e.* point  $v_a$  and point  $v_b$ ) are detected from distinct classes of the object during optimization. The collision loss shrinks the object size such that there are no different clusters of negative SDF points on the object surface can be detected. (c). The *stability maintenance* keeps the the child and the parent to be as close as possible.

**Physical Plausibility.** Physical plausibility ensure generated 3D scenes adhering to fundamental physical principles. Instead of using 3D bounding box (3DBB) as object approximation, *Scenethesis* accurately detects collision state from the surface points of the 3D models using signed distance field (SDF). The *collision avoidance* and *stability maintenance* as illustrated in Figure 10.

The *collision avoidance* affects the translation  $T$  by:

$$\mathcal{L}_{\text{translation}} = \|f(\mathbf{T}, |d|, \mathbf{u}) - \mathbf{T}\|_2^2, \quad (6)$$

where  $f(\mathbf{T}, d, \mathbf{u}) = \mathbf{T} + \mathbf{u} \cdot |d|$  computes a collision-free position  $\hat{\mathbf{T}}$  by adjusting  $\mathbf{T}$  along direction  $\mathbf{u}$  with step size  $d$ . Here,  $d$  is the negative SDF value at a collided point  $v_i$  such that  $d(v_i) \leq 0$  and  $|d| = \max(0, -d(v_i))$  is the negative SDF value  $d$  after being processed through a ReLU function, meaning only collided points contribute to this collision term. The direction  $\mathbf{u}$  is defined from the collision point toward the model’s centroid  $\mathbf{C}$ , guiding objects away from collisions.

The *collision avoidance* affects the scaling  $s$  by detecting that object collides from at least two different directions:

$$\mathcal{L}_{\text{scale}} = \begin{cases} \sum_{v_i \in \mathbf{V}^-} \left( g(|d_i|, \mathbf{u}_i) - s \right)^2 & \text{if } N_{\text{cluster}} > 1 \\ 0 & \text{otherwise} \end{cases}, \quad (7)$$

where  $g(|d_i|, \mathbf{u}_i) = \frac{||\mathbf{u}_i|| - |d_i|}{||\mathbf{u}_i||}$  defines the target scale to reduce collision regions.  $N_{\text{cluster}}$  denotes the number of distinct clusters formed without SDF sign flipping. As shown in Figure 4, two surface points  $i$  and  $j$  with  $d_i \leq 0$  and  $d_j \leq 0$  belong to different clusters, and thus push the object to be smaller.

The *stability maintenance* affects the translation  $T$  by:

$$\mathcal{L}_{\text{stability}} = \sum_{v_i \in \mathbf{V}^B} \left( 1 - \exp(-d_i^2) \right), \quad (8)$$

where  $\mathbf{V}^B$  are the sampled points on the bottom surface of bounding box, and  $d_i$  are their corresponding SDF values.

**Method Overview** Building on the physics-aware optimization described above, we now integrate pose spatial constraints and physical constraints into the text-to-3D optimization framework. Since physical loss depends on object geometry and can alter its spatial position, it may affect the rendered visible regions due to occlusions or shifts in the scene, introducing biases in pose alignment when using image guidance for semantic correspondence matching. To mitigate this, we adopt a two-stage optimization strategy: first, we optimize pose alignment based on correspondence matching; then, we refine object placement with physical constraints to ensure a visually coherent and physically plausible 3D scene. The following function defines the joint optimization of object position, orientation, and scale:

$$\mathcal{L} = \lambda_p \mathcal{L}_{\text{pose}} + \lambda_{c,T} \mathcal{L}_{\text{translation}} + \lambda_{c,S} \mathcal{L}_{\text{scale}} + \lambda_s \mathcal{L}_{\text{stability}} \quad (9)$$

### 6.3. Experiment Details

**Parameters.** For pose alignment, we select the  $m = 100$  correspondence points with matching conference  $\tau \geq 0.6$  in each optimization iteration. Additionally, we uniformly select  $n = 400$  samples from the surface of 3D model to



Figure 11. *Short prompt*: a living room with reading materials; *detailed long prompt*: A living room that provide a neutral and cozy space with a minimalist design. At the center of the scene, a light beige sofa is positioned against a textured stone wall in the background. In front of the sofa, a round wooden coffee table sits on the floor, with a white coffee cup placed on top. Two blue armchairs are symmetrically arranged on either side of the coffee table, facing inward toward the sofa. Behind each armchair, a tall white floor lamp stands, providing ambient lighting. Next to the lamps, green potted plants are placed near the wall, adding a natural decorative touch.

accurately detect the collision and stability states in each optimization iteration.

We explored *Adam* and *SGD* as the optimizer during the optimization process. Though *Adam* has been widely applied for training deep neural networks, the adaptive momentum makes the optimization unstable, leading to sub-optimal optimized pose. So we use *SGD* in our implementation. The optimization implementation is based on pytorch3D [38] and the visualization is rendered using Blender.

**Prompts.** *Scenethesis* supports both short and detailed user-specified prompts. A short prompt provides a user-friendly and flexible approach to 3D scene generation, where the LLM interprets the input, revisits the available 3D models in database, selects the common objects and anchor objects, and generates an upsampled text prompt for coarse layout planning. In contrast, a long prompt, which includes user-defined objects and inter-object relationships, enables greater user control over 3D scene generation. In this case, the LLM directly reasons over the detailed prompt, revisits available 3D models in the database, and defines the anchor object, skipping the upsampling stage. We illustrate examples of short and long prompts defining a living room in Figure 11.

We compared four baselines—Physcene [47], Diffuscene [41], SceneTeller [33], and Holodeck [48]—evaluating visual quality, physical plausibility, and interactivity metrics. Among them, Diffuscene, SceneTeller, and Holodeck perform text-to-3D scene generation. For visual quality assessment, we use both a user study and GPT-4o as evaluation tools. Unlike other baselines, which generate only living room, bedroom,

	A is better	B is better
Layout coherence	<input type="radio"/>	<input type="radio"/>
Spatial realism	<input type="radio"/>	<input type="radio"/>
Overall consistency	<input type="radio"/>	<input type="radio"/>

Figure 12. User study example.

and dining room scenes from 3D-FRONT [12], Holodeck and *Scenethesis* utilize Objaverse [8] as a retrieval database, enabling more diverse indoor scene generation.

We outline the GPT-4o prompt assessment for both baseline evaluation and ablation evaluation as follows:

- **Comparison with baselines by GPT-4o:** GPT-4o is employed to evaluate the generated scenes for four metrics: *object diversity*, *layout coherence*, *spatial realism* and *complexity*, and *overall performance*. The evaluation prompts are detailed in the *Instruction Prompts for Evaluating Generated Scene* Section 7.3. Additionally, a comparison example of the generated scenes is provided in Figure 13 with their evaluation results generated by GPT-4o detailed in *Evaluation Example of Generated Scenes* Section 7.4.
- **Comparison with baselines by human preference:** We applied a user study to study human preference of baseline method with our method. See Figure 12 as an example. There are 69 users took our survey.
- **Evaluation in Ablation Studies:** GPT-4o is also utilized to assess the pose alignment metric during the ablation studies of *Scenethesis*'s physics-aware optimization. This evaluation measures the similarity of object position, size, and orientation with their counterparts in the image guidance, as well as the overall visual coherence of the layout. The instructions for assessing pose alignment are provided in the *Instruction Prompts for Ablation Study* Section 7.5.

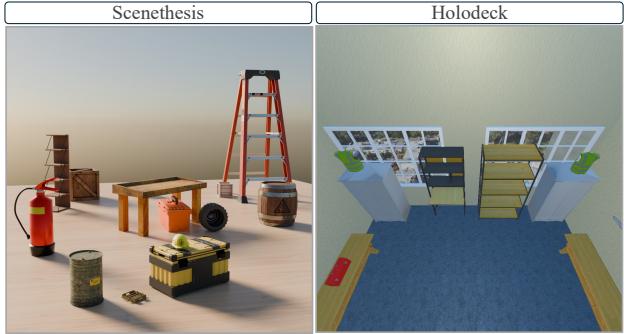


Figure 13. An example comparison of generated scenes given user prompt: “a warehouse”. Note that *Scenethesis*’s scenes are rendered without an environment map to ensure a fair comparison with Holodeck’s scenes.

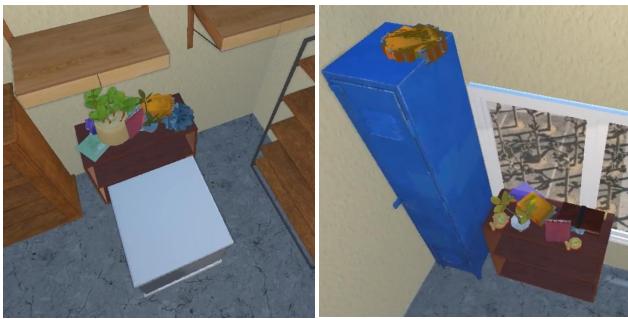


Figure 14. An example of objects collision from Holodeck’s scenes.

**Results.** We present additional qualitative results for *Scenethesis*’s scenes.

- **Qualitative Results of *Scenethesis*’s Scene:** We present different camera views to showcase the qualitative results of *Scenethesis*’s scenes, as shown in Figure 15. Figure 17 presents the generated scenes by *Scenethesis* with their image guidance. Note that *Scenethesis* focuses on layout planning for ground objects. The absence of certain unique assets in Objaverse [8] may cause discrepancies between the generated scene and the image guidance. Future work could address this by incorporating more diverse assets.
- **Quantitative Results of Physical plausibility Comparison:** The physical Plausibility quantitative comparison presented in Table 1 of the **Experiment** section. While Holodeck applies both soft and hard constraints based on the Depth-First-Search Solver and small objects are placed on predefined locations. These small objects may collide with each other due to the shape and size variations as shown in Figure 14.
- **Visual Comparison with Holodeck:** In addition to the quantitative comparison presented in Table 1 of the **Experiment** section, we provide a visual comparison between *Scenethesis* and Holodeck, a state-of-the-art LLM-based 3D interactive scene generation method, in Figure 16. Based on the four evaluation metrics detailed in the **Ex-**

**periment** section, scenes generated by *Scenethesis* demonstrate greater diversity in object categories, quantities, and sizes. More importantly, *Scenethesis*’s scenes have a broader range of spatial relationships, such as “*on top of*”, “*inside*”, and “*under*”, compared to those generated by Holodeck [48], which supports only “*on top of*” spatial relation. Furthermore, *Scenethesis*’s scenes align more faithfully with the intended scene type. *i.e.* when given the description “*a peaceful beach during sunset*”, *Scenethesis* produces an outdoor scene with appropriate beach elements, while Holodeck incorporates beach-related objects but generates an environment resembling an indoor setting.



Figure 15. Qualitative results of generated indoor and outdoor scenes by *Scenethesis* at different camera viewpoints

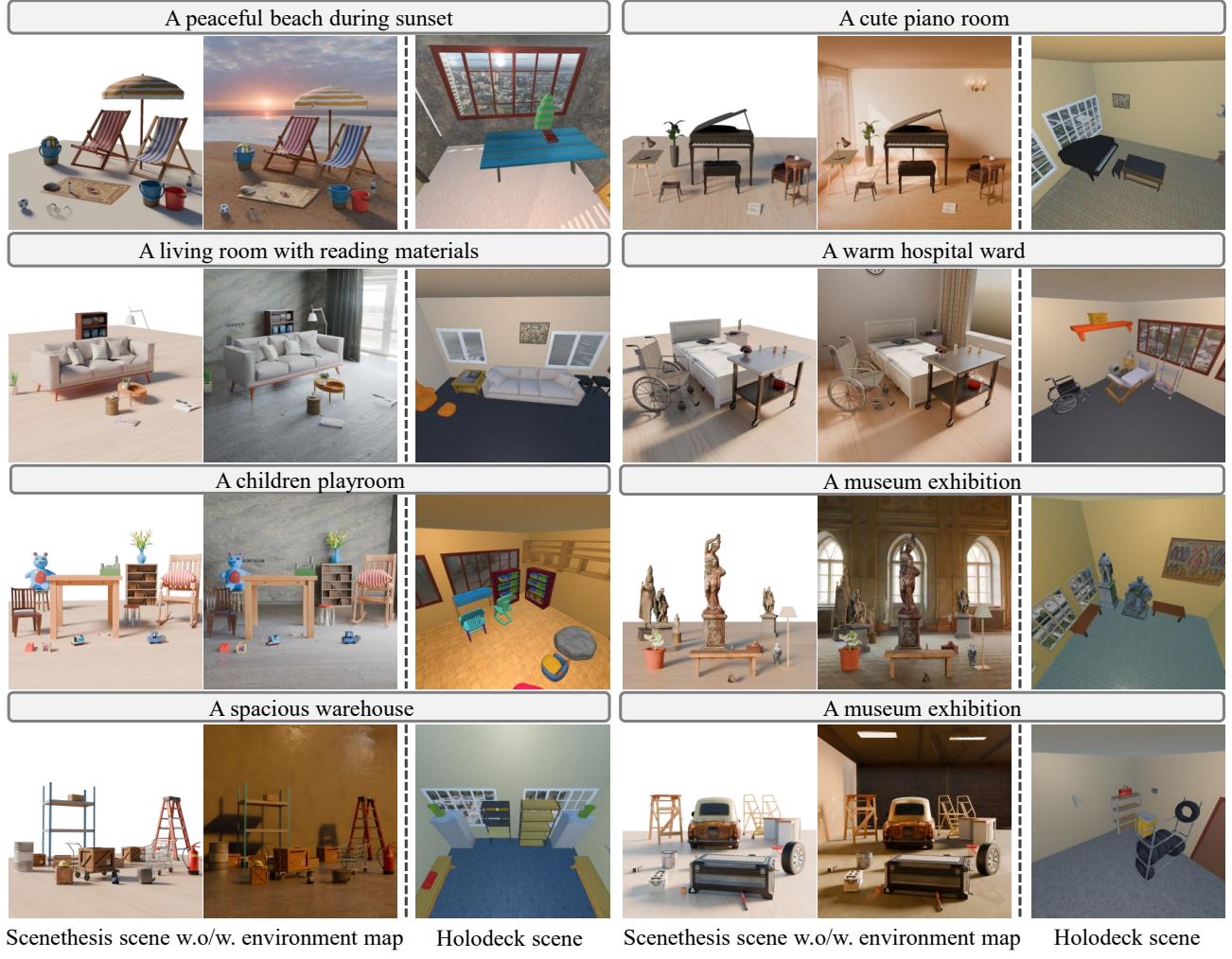


Figure 16. Visualization comparison of generated scenes between *Scenethesis* and Holodeck. The first column of images shows scenes generated by *Scenethesis* without an environment map, the second column displays scenes generated by *Scenethesis* with an environment map, and the third column presents scenes generated by Holodeck. The evaluation metrics, including *object diversity*, *layout coherence*, *spatial realism*, and *overall performance*, are detailed in the **Experiment** section. *Scenethesis*'s scenes have a wider variety of spatial relationships, such as “*on top of*”, “*inside*”, and “*under*”, compared to those generated by Holodeck [48], which supports only “*on top of*” spatial relation. In addition, Holodeck lacks visual perception and usually generates misoriented objects, e.g. shelves occlude the window in children playroom and warehouse, chair orients towards the window in the hospital case, hindering their functionalities.

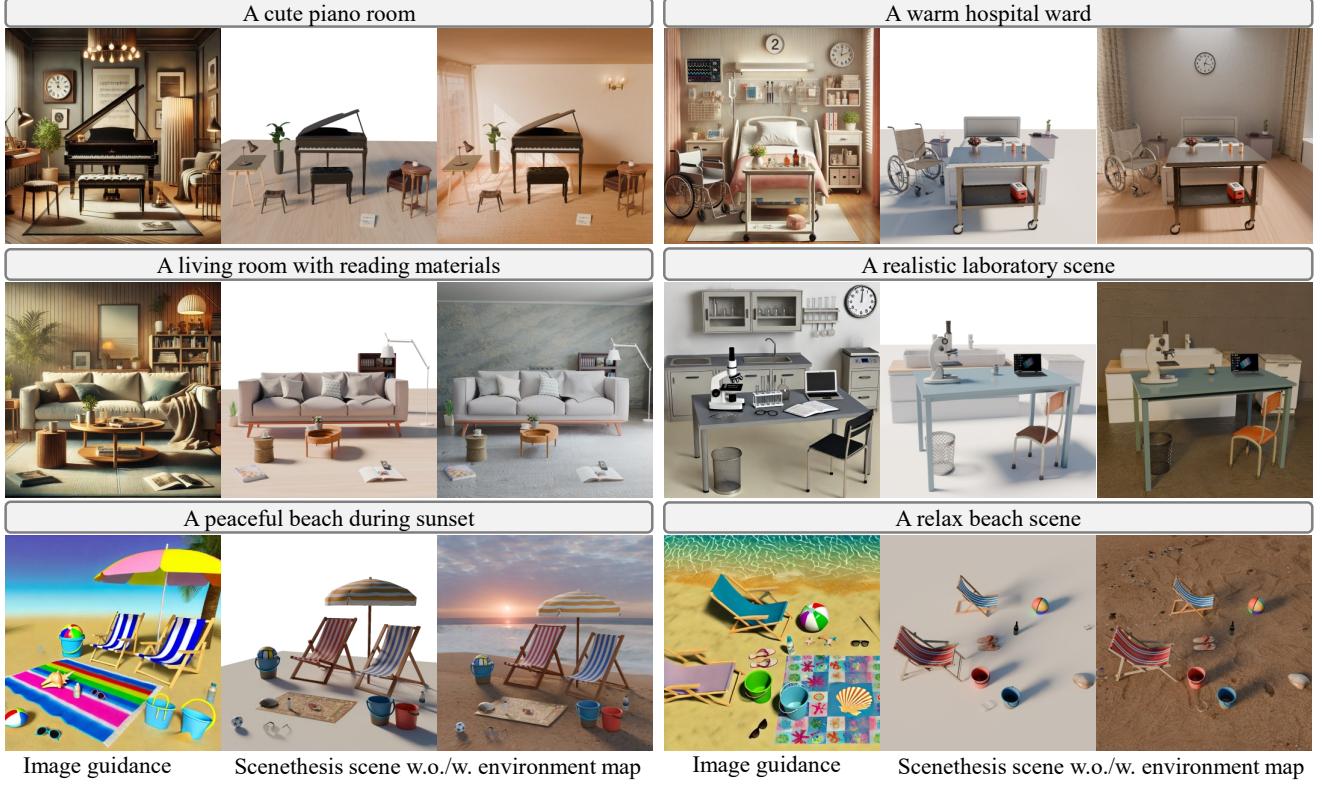


Figure 17. We provide a visual illustration of the generated scenes and their corresponding image guidance. The first column displays the image guidance, while the second and third columns show the generated scenes without and with the environment map, respectively. Note that *Scenethesis* focus on layout planning for objects on the ground. Additionally, certain unique assets, such as a beach mat, are unavailable in Objaverse [8], which may result in the generated scene differing from the image guidance. Future work could enhance the system by incorporating a wider range of assets.

## 7. Prompts Examples

### 7.1. Coarse Scene Planning Instruction Prompts

#### Coarse Scene Planning Instruction Prompts

##### Task Description:

You are responsible for generating a set of common objects and planning a scene based on these common objects. You will be given a list that includes all available object categories and a text prompt to describe a scene. This is a hard task, please think deeply and write down your analysis in following steps:

##### Step 1: Review All Categories

- Begin by thoroughly reviewing the categories in the provided list.
- Identify potential groups or clusters of objects within this list that are commonly found in similar environments (e.g., furniture, electronics, household items, etc.).

##### Step 2: Interpret Input Prompt

- Carefully read the input prompt. Understand the theme, primary activities, or the setting it describes, as these will guide your object selection. *i.e.* if the prompt gives: *children playing room*, then you may think of objects like tent, toy, bear, ball, chair, etc.

##### Step 3: Object Selection

- Based on the description, select at least 15 object categories from the list that match the scene.
- Determine the anchor object:
  - Identify the anchor object among the selected objects. Consider the following factors:
    - A large object directly on the ground (*i.e.* floor, table, or shelf).
    - An object that influences where other objects are placed (*i.e.* a table in a dining room, and there are cups and fruits on the table).
    - The object should logically anchor the scene and often defines the scene's layout orientation. *i.e.* the sofa in a front-facing view in the scene.

##### Step 4: Object Cross-check

- I will give you \$100 tips if you can cross-check whether objects in the scene can be found in the given category list or its relevant categories. *i.e.*, if there is a bookshelf in your planned scene, the bookshelf should also be found in the given list, or bookcase can be found in the list if bookshelf is not covered by the category. Otherwise, re-plan the scene.

##### Step 5: Plan Scene with Selected Objects

- Based on the description and selected objects, plan the scene, keeping these aspects in mind:
  - Functionality:** Choose objects that are contextually relevant to the scene (e.g., selecting a table, chair, flower vase, and utensils for a dining room), but do not generate any wall décor objects.
  - Spatial Hierarchy:**
    - Please have a depth effect in the layout. For the depth effect, the scene should have some objects placed on the ground as the background, central, and in the front, resulting in a depth layout. *i.e.* the sofa and bookshelf are the background of the table and chair set in the living room.
    - Please have a supportive item in the layout. *i.e.* the shoes, bag, and hat are in the display shelf in a clothes store, where the display shelf is a supportive item.
  - Balance:** Ensure a mix of large and small objects to avoid overcrowding or under-populating the scene. *i.e.* taking the table as the center, there are flower vases, fruits, and cups on the table, and chairs are on the sides.

##### Step 6: Output Format:

- Save the selected objects as a json file follow the output format:

**Anchor object:**

**Other common objects:**

- Save scene planning as txt file.

## 7.2. Coarse Scene Planning Output Example

### Coarse Scene Planning Output Example

#### User input: a gym room

##### Step 1: Reviewing Object Categories for Gym-Related Items

From the provided list, here are the relevant categories typically found in a gym setting:

- Exercise equipment: dumbbell, treadmill, weight, barbell
- Furniture or storage items: bench, shelf, cabinet
- Accessories and smaller items: towel, water bottle, yoga mat, scale, bag

##### Step 2: Scene Interpretation

The scene is a gym room, so the primary objects would focus on equipment for strength and cardio exercises, along with supportive items like shelves for organization and accessories like towels and water bottles for hydration.

##### Step 3: Scene Layout and Object Selection

- Anchor object: Treadmill (as the main focus of cardio equipment)
- Other common objects: Dumbbell, Barbell, Bench, Shelf, Water bottle, Towel, Weight, Bag, Mat, Scale, Cabinet, Exercise ball, Fan, Mirror

##### Step 4: Object cross-check

- Treadmill - ✓ Available as "treadmill"
- Dumbbell - ✓ Available as "dumbbell"
- Barbell - ✓ Available as "barbell"
- Bench - ✓ Available as "bench"
- Shelf - ✓ Available as "shelf"
- Water bottle - ✓ Available as "water bottle"
- Towel - ✓ Available as "towel rack"
- Weight - ✓ Available as "weight"
- Bag - ✓ Available as "bag"
- Mat - Not explicitly listed, but a close alternative like "place mat"
- Scale - ✓ Available as "scale (measuring instrument)"
- Cabinet - ✓ Available as "cabinet"
- Exercise ball - Not explicitly listed, but a close alternative like "ball"
- Fan - ✓ Available as "fan"
- Mirror - ✓ Available as "mirror"

##### Step 5: Plan Scene with Selected Objects

- Background: Shelf and cabinet are placed to the side or back corner of the room for organization, holding smaller items. Mirror mounted on the back wall for checking exercise form. They serve as supportive equipments...It enhances depth and makes the room feel more spacious...
- Middle Ground: Treadmill is the primary cardio equipment and anchor object, anchoring the scene. Other objects such as bench and barbell are placed around the treadmill, allowing access from multiple directions for various exercises. They provide visual weight to the scene and represent strength-focused equipment...
- Foreground: Dumbbell, water bottle and other objects set closer to the viewer, easily accessible for hand-held exercises. It introduces a smaller, more detailed object in the front. They placed nearby, adding another layer of exercise options for core stability...

This hierarchy ensures a balanced layout, where larger and bulkier items form the background... The middle ground includes the treadmill and training equipment, while smaller and accessible items like dumbbells occupy the foreground... This setup provides a clear layout for practical use and visual depth in the gym scene...

##### Step 6: Output Format

A json file that include the available objects and a txt file that describe the scene planning.

### 7.3. Instruction for Evaluating Generated Scenes

#### GPT-4o Evaluation Prompts for Generated Scenes

##### Evaluate Generated Scenes for *Scenethesis* and Holodeck

This is a hard problem. You are supposed to compare the alignment of a pair of images with a given text prompt that describe the scene. Images contain generated scenes by two different methods. Please evaluate them in the following five aspects:

1. **Object diversity:** Counting number of object and object types in the scene. The higher number, the better object diversity. *i.e.* if there are three shelves and a box in the scene. Then the object type is 2 and the number of counting is  $3+1 = 4$ .
2. **Layout coherence:** whether the objects position and orientation in the scene are realistic and adhere to common sense.
  - An ideal layout would be dependents on the scene type. *i.e.* the objects in the garage scene will forms as a bit chaotic organization and it would be less reasonable if all objects standing against the wall in a very clean order.
  - Objects should be placed reasonably. *i.e.* shelves fly in the air or hange on the wall.
3. **Spatial realism & complexity:** whether scene contains diverse hierarchy. The measure of spatial complexity is by review scenes and carefully evaluate objects relations. The higher diverse relation indicate a better spatial complexity. The spatial hierarchy refers to the following aspects:
  - relations such as on the top of, in side of, under the, etc.
  - objects in various sizes. *i.e.* the scene should contain large objects such as shelf and small objects such as apple.
  - objects that are arranged in different area and form a depth effects. *i.e.* objects are in the background, in the front, and in the center of the space.
4. **Overall performance:** assets' categorizes and styles are faithful to the scene type. *i.e.* fruits are less likely in the warehouse, while carts, boxes are more likely in a warehouse scene. Thus, it would be a bad performance if find fruits in the warehouse and good performance if cartsshelfboxes presents in the warehouse. Please carefully evaluate the object's category and scene type. To see whether it is reasonable enough based on your common sense.
5. Please ignore aspects in : environment background style such as lighting, shadow, color, and decorations in the wall etc. This is a challenging task, please think aloud and make your decisions by:
  - analyzing each criteria one by one for each image.
  - report your evaluation results based on the analysis.
6. Your *output format* should be: {} is better than {} in the criteria {}.

## 7.4. Evaluation Example of Generated Scenes

### GPT-4o Evaluation Results

#### Scenethesis's Scene

##### 1. Object diversity:

- Objects visible: Ladder, barrel, boxes, fire extinguisher, table, tire, crate.
- Object types: 8 types. Total count: 13
- Higher diversity due to various objects.

##### 2. Layout coherence:

- Objects are dispersed naturally, suggesting activity.
- Different sizes and a less structured arrangement suggest a more realistic warehouse scene.

##### 3. Spatial realism & complexity:

- Objects have various spatial relations; ladder stands, boxes on ground, tables at center...
- Good diversity with a mix of background and central objects.

##### 4. Overall performance:

- Varied objects like ladders, barrels, crates, and extinguishers are typical in a warehouse.
- Realistic depiction of a warehouse setting.

#### Holodeck's Scene

##### 1. Object diversity:

- Objects visible: Shelves, cabinets, safety vests, and workbenches.
- Object types: Shelves, cabinets, vests, workbenches. Total count:  $7+2=9$  (shelves are considered per structure, as are others).
- Diversity is somewhat low as there are few types.

##### 2. Layout coherence:

- Organization is neat; workbenches and cabinets are along the walls.
- Reasonable layout for a tidy scene, but lacks variety in sizes.

##### 3. Spatial realism & complexity:

- Shelving and cabinets against walls create vertical space; workbenches on the ground.
- Lacks diverse spatial relations.

##### 4. Overall performance:

Objects like shelves and workbenches fit a warehouse setting, but it lacks variety typical in larger-scale warehousing.

#### Evaluation

**Object diversity:** *Scenethesis* is better than Holodeck; **Layout coherent:** *Scenethesis* is better than Holodeck;

**Spatial realism & complexity:** *Scenethesis* is better than Holodeck; **Overall performance:** *Scenethesis* is better than Holodeck.

## 7.5. Instruction Prompts for Ablation Study

### Pose Alignment Evaluation Instruction Prompt

This task involves evaluating the pose alignment between two images in a pair. One image serves as the image guidance (GT), while the other is a generated image. Your objective is to measure the pose alignment of the generated image relative to the GT image. Follow these steps for evaluation:

1. **Review Objects in the GT Image:** Examine all objects in the GT image, focusing on their locations, sizes, and orientations. Understand the spatial relationships among objects, such as *on top of*, *inside*, *under*, etc.
2. **Evaluate pose alignment:** Assess the similarity between the generated image and the GT image based on the following three aspects:
  - Location and Size Similarity: Compare the location and size of objects in the generated image with those in the GT image. Assign a similarity score between 0 and 1, where 1 indicates the highest similarity. For example:
    - If an apple in the GT image is placed at the center of a table, and in the generated image it is placed on the left side of the table, the similarity might be moderate (e.g., 0.5).
    - If the apple is misplaced (e.g., on the ground or missing entirely), the similarity would be very low (e.g., 0.1).
  - Orientation Similarity: Examine the orientation of each object in the generated image compared to the GT image. Pay close attention to details, noting any deviations such as slight tilts (e.g., right/left, up/down) or rotations that create different perspectives. Assign a score from 0 to 1, where 1 indicates perfect alignment and 0 indicates a significant mismatch (e.g., opposite orientation).
  - Overall Layout Similarity: Assess the overall visual coherence of the generated image compared to the GT image, including spatial relationships and hierarchical structure. Assign a similarity score between 0 and 1, where 1 represents a perfect match. For instance:
    - A perfect match occurs when the generated image maintains the same spatial relationships, relative locations, sizes, and orientations as the GT image (e.g., an apple placed at the center of a table in both images).
    - Small deviations in placement or orientation are acceptable but should result in a lower score.
3. **Exclusions:** Do not consider style, appearance, object shape, or texture in your evaluation. Focus solely on pose alignment.
4. **Output Format:** Clearly document your similarity scores for each aspect (Location and Size Similarity, Orientation Similarity, and Overall Layout Similarity) following the format: location and size similarity score is {}, orientation similarity score is {}, and overall layout similarity score is {}. Please save the evaluated scores as a json file.

## References

- [1] Rio Aguina-Kang, Maxim Gumin, Do Heon Han, Stewart Morris, Seung Jean Yoo, Aditya Ganeshan, R Kenny Jones, QiuHong Anna Wei, Kailiang Fu, and Daniel Ritchie. Open-universe indoor scene generation using llm program synthesis and uncurated object databases. *arXiv preprint arXiv:2403.09675*, 2024. 3
- [2] Aleksei Bochkovskii, Amaël Delaunoy, Hugo Germain, Marcel Santos, Yichao Zhou, Stephan R Richter, and Vladlen Koltun. Depth pro: Sharp monocular metric depth in less than a second. *arXiv preprint arXiv:2410.02073*, 2024. 4, 1
- [3] Angel Chang, Manolis Savva, and Christopher D Manning. Learning spatial knowledge for text to 3d scene generation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 2028–2038, 2014. 3
- [4] Yamei Chen, Yan Di, Guangyao Zhai, Fabian Manhardt, Chenyangguang Zhang, Ruida Zhang, Federico Tombari, Nassir Navab, and Benjamin Busam. Secondpose: Se (3)-consistent dual-stream feature fusion for category-level pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9959–9969, 2024. 4
- [5] Yunuo Chen, Tianyi Xie, Zeshun Zong, Xuan Li, Feng Gao, Yin Yang, Ying Nian Wu, and Chenfanfu Jiang. Atlas3d: Physically constrained self-supporting text-to-3d for simulation and fabrication. *arXiv preprint arXiv:2405.18515*, 2024. 6
- [6] Matt Deitke, Eli VanderBilt, Alvaro Herrasti, Luca Weihs, Kiana Ehsani, Jordi Salvador, Winson Han, Eric Kolve, Aniruddha Kembhavi, and Roozbeh Mottaghi. Proctor: Large-scale embodied ai using procedural generation. *Advances in Neural Information Processing Systems*, 35:5982–5994, 2022. 2, 3
- [7] Matt Deitke, Rose Hendrix, Ali Farhadi, Kiana Ehsani, and Aniruddha Kembhavi. Phone2proc: Bringing robust robots into our chaotic world. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9665–9675, 2023. 2, 3
- [8] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13142–13153, 2023. 4, 6, 5, 8
- [9] Johan Edstedt, Qiyu Sun, Georg Bökman, Mårten Wadenbäck, and Michael Felsberg. Roma: Robust dense feature matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19790–19800, 2024. 4, 5, 1
- [10] Weixi Feng, Wanrong Zhu, Tsu-jui Fu, Varun Jampani, Arjun Akula, Xuehai He, Sugato Basu, Xin Eric Wang, and William Yang Wang. Layoutgpt: Compositional visual planning and generation with large language models. *Advances in Neural Information Processing Systems*, 36, 2024. 2
- [11] Rafail Friedman, Amit Abecasis, Yoni Kasten, and Tali Dekel. Scenescape: Text-driven consistent scene generation. *Advances in Neural Information Processing Systems*, 36, 2024. 3
- [12] Huan Fu, Rongfei Jia, Lin Gao, Mingming Gong, Binqiang Zhao, Steve Maybank, and Dacheng Tao. 3d-future: 3d furniture shape with texture. *International Journal of Computer Vision*, 129:3313–3337, 2021. 2, 3, 6, 7, 4
- [13] Chuang Gan, Jeremy Schwartz, Seth Alter, Damian Mrowca, Martin Schrimpf, James Traer, Julian De Freitas, Jonas Kubilius, Abhishek Bhandwaldar, Nick Haber, et al. Threedworld: A platform for interactive multi-modal physical simulation. *arXiv preprint arXiv:2007.04954*, 2020. 2
- [14] Lukas Höller, Ang Cao, Andrew Owens, Justin Johnson, and Matthias Nießner. Text2room: Extracting textured 3d meshes from 2d text-to-image models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7909–7920, 2023. 2, 3
- [15] Ziniu Hu, Ahmet Iscen, Aashi Jain, Thomas Kipf, Yisong Yue, David A Ross, Cordelia Schmid, and Alireza Fathi. Scenecraft: An llm agent for synthesizing 3d scenes as blender code. In *Forty-first International Conference on Machine Learning*, 2024. 2
- [16] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024. 4, 6, 1
- [17] Mukul Khanna, Yongsen Mao, Hanxiao Jiang, Sanjay Haresh, Brennan Shacklett, Dhruv Batra, Alexander Clegg, Eric Undersander, Angel X Chang, and Manolis Savva. Habitat synthetic scenes dataset (hssd-200): An analysis of 3d scene scale and realism tradeoffs for objectgoal navigation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16384–16393, 2024. 3
- [18] Eric Kolve, Roozbeh Mottaghi, Winson Han, Eli VanderBilt, Luca Weihs, Alvaro Herrasti, Matt Deitke, Kiana Ehsani, Daniel Gordon, Yuke Zhu, et al. Ai2-thor: An interactive 3d environment for visual ai. *arXiv preprint arXiv:1712.05474*, 2017. 2
- [19] Jacob Krantz, Erik Wijmans, Arjun Majumdar, Dhruv Batra, and Stefan Lee. Beyond the nav-graph: Vision-and-language navigation in continuous environments. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXVIII 16*, pages 104–120. Springer, 2020. 2
- [20] Vikram Kumaran, Jonathan Rowe, Bradford Mott, and James Lester. Scenecraft: Automating interactive narrative scene generation in digital games with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment*, pages 86–96, 2023. 2, 3
- [21] Chengshu Li, Ruohan Zhang, Josiah Wong, Cem Gokmen, Sanjana Srivastava, Roberto Martín-Martín, Chen Wang, Gabrael Levine, Michael Lingelbach, Jiankai Sun, et al. Behavior-1k: A benchmark for embodied ai with 1,000 everyday activities and realistic simulation. In *Conference on Robot Learning*, pages 80–93. PMLR, 2023. 2
- [22] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen

- image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR, 2023. 6
- [23] Manyi Li, Akshay Gadi Patil, Kai Xu, Siddhartha Chaudhuri, Owais Khan, Ariel Shamir, Changhe Tu, Baoquan Chen, Daniel Cohen-Or, and Hao Zhang. Grains: Generative recursive autoencoders for indoor scenes. *ACM Transactions on Graphics (TOG)*, 38(2):1–16, 2019. 3
- [24] Yixun Liang, Xin Yang, Jiantao Lin, Haodong Li, Xiaogang Xu, and Yingcong Chen. 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*, pages 6517–6526, 2024. 4
- [25] Chenguo Lin and Yadong Mu. Instructscene: Instruction-driven 3d indoor scene synthesis with semantic graph prior. *arXiv preprint arXiv:2402.04717*, 2024. 3
- [26] Tsung-Yi Lin, Chen-Hsuan Lin, Yin Cui, Yunhao Ge, Seungjun Nah, Arun Mallya, Zekun Hao, Yifan Ding, Hanzi Mao, Zhaozhou Li, et al. Genusd: 3d scene generation made easy. In *ACM SIGGRAPH 2024 Real-Time Live!*, pages 1–2. 2024. 3
- [27] Zhiqiu Lin, Deepak Pathak, Baiqi Li, Jiayao Li, Xide Xia, Graham Neubig, Pengchuan Zhang, and Deva Ramanan. Evaluating text-to-visual generation with image-to-text generation. In *European Conference on Computer Vision*, pages 366–384. Springer, 2024. 6
- [28] Zhiqiu Lin, Deepak Pathak, Baiqi Li, Jiayao Li, Xide Xia, Graham Neubig, Pengchuan Zhang, and Deva Ramanan. Evaluating text-to-visual generation with image-to-text generation. In *European Conference on Computer Vision*, pages 366–384. Springer, 2025. 6
- [29] Lu Ling, Yichen Sheng, Zhi Tu, Wentian Zhao, Cheng Xin, Kun Wan, Lantao Yu, Qianyu Guo, Zixun Yu, Yawen Lu, et al. Dl3dv-10k: A large-scale scene dataset for deep learning-based 3d vision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22160–22169, 2024. 6
- [30] Andrew Luo, Zhoutong Zhang, Jiajun Wu, and Joshua B Tenenbaum. End-to-end optimization of scene layout. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3754–3763, 2020. 3
- [31] Miles Macklin. Warp: A high-performance python framework for gpu simulation and graphics. <https://github.com/nvidia/warp>, 2022. NVIDIA GPU Technology Conference (GTC). 6
- [32] Soroush Nasiriany, Abhiram Maddukuri, Lance Zhang, Adeet Parikh, Aaron Lo, Abhishek Joshi, Ajay Mandlekar, and Yuke Zhu. Robocasa: Large-scale simulation of everyday tasks for generalist robots. *arXiv preprint arXiv:2406.02523*, 2024. 2, 8
- [33] Başak Melis Öcal, Maxim Tatarchenko, Sezer Karaoglu, and Theo Gevers. Sceneteller: Language-to-3d scene generation. In *European Conference on Computer Vision*, pages 362–378. Springer, 2024. 2, 3, 6, 7, 4
- [34] Despoina Paschalidou, Amlan Kar, Maria Shugrina, Karsten Kreis, Andreas Geiger, and Sanja Fidler. Atiss: Autoregressive transformers for indoor scene synthesis. *Advances in Neural Information Processing Systems*, 34:12013–12026, 2021. 2, 3
- [35] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019. 6
- [36] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 6
- [37] Alexander Raistrick, Lingjie Mei, Karhan Kayan, David Yan, Yiming Zuo, Beining Han, Hongyu Wen, Meenal Parakh, Stamatis Alexandropoulos, Lahav Lipson, et al. Infinigen indoors: Photorealistic indoor scenes using procedural generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21783–21794, 2024. 3
- [38] Nikhila Ravi, Jeremy Reizenstein, David Novotny, Taylor Gordon, Wan-Yen Lo, Justin Johnson, and Georgia Gkioxari. Accelerating 3d deep learning with pytorch3d. *arXiv preprint arXiv:2007.08501*, 2020. 6, 4
- [39] Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kunchang Li, He Cao, Jiayu Chen, Xinyu Huang, Yukang Chen, Feng Yan, et al. Grounded sam: Assembling open-world models for diverse visual tasks. *arXiv preprint arXiv:2401.14159*, 2024. 4, 1
- [40] Fan-Yun Sun, Weiyu Liu, Siyi Gu, Dylan Lim, Goutam Bhat, Federico Tombari, Manling Li, Nick Haber, and Jiajun Wu. Layoutlm: Differentiable optimization of 3d layout via vision-language models. *arXiv preprint arXiv:2412.02193*, 2024. 3
- [41] Jiapeng Tang, Yinyu Nie, Lev Markhasin, Angela Dai, Justus Thies, and Matthias Nießner. Diffuscene: Denoising diffusion models for generative indoor scene synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 20507–20518, 2024. 2, 3, 6, 7, 4
- [42] Kai Wang, Yu-An Lin, Ben Weissmann, Manolis Savva, Angel X Chang, and Daniel Ritchie. Planit: Planning and instantiating indoor scenes with relation graph and spatial prior networks. *ACM Transactions on Graphics (TOG)*, 38(4):1–15, 2019. 3
- [43] Xinpeng Wang, Chandan Yeshwanth, and Matthias Nießner. Sceneformer: Indoor scene generation with transformers. In *2021 International Conference on 3D Vision (3DV)*, pages 106–115. IEEE, 2021. 3
- [44] Yufei Wang, Zhou Xian, Feng Chen, Tsun-Hsuan Wang, Yian Wang, Katerina Fragkiadaki, Zackory Erickson, David Held, and Chuang Gan. Robogen: Towards unleashing infinite data for automated robot learning via generative simulation. *arXiv preprint arXiv:2311.01455*, 2023. 3
- [45] Yian Wang, Xiaowen Qiu, Jiageng Liu, Zhehuan Chen, Jiting Cai, Yufei Wang, Tsun-Hsuan Johnson Wang, Zhou Xian, and Chuang Gan. Architect: Generating vivid and interactive 3d

- scenes with hierarchical 2d inpainting. *Advances in Neural Information Processing Systems*, 37:67575–67603, 2024. 3
- [46] Rundi Wu, Ben Mildenhall, Philipp Henzler, Keunhong Park, Ruiqi Gao, Daniel Watson, Pratul P Srinivasan, Dor Verbin, Jonathan T Barron, Ben Poole, et al. Reconfusion: 3d reconstruction with diffusion priors. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 21551–21561, 2024. 4
- [47] Yandan Yang, Baoxiong Jia, Peiyuan Zhi, and Siyuan Huang. Physcene: Physically interactable 3d scene synthesis for embodied ai. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16262–16272, 2024. 2, 3, 5, 6, 7, 8, 4
- [48] Yue Yang, Fan-Yun Sun, Luca Weihs, Eli VanderBilt, Alvaro Herrasti, Winson Han, Jiajun Wu, Nick Haber, Ranjay Krishna, Lingjie Liu, et al. Holodeck: Language guided generation of 3d embodied ai environments. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16227–16237, 2024. 2, 3, 4, 5, 6, 7
- [49] Hong-Xing Yu, Haoyi Duan, Charles Herrmann, William T Freeman, and Jiajun Wu. Wonderworld: Interactive 3d scene generation from a single image. *arXiv preprint arXiv:2406.09394*, 2024. 2, 3
- [50] Hong-Xing Yu, Haoyi Duan, Junhwa Hur, Kyle Sargent, Michael Rubinstein, William T Freeman, Forrester Cole, Deqing Sun, Noah Snavely, Jiajun Wu, et al. Wonderjourney: Going from anywhere to everywhere. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6658–6667, 2024. 3
- [51] Junyi Zhang, Charles Herrmann, Junhwa Hur, Eric Chen, Varun Jampani, Deqing Sun, and Ming-Hsuan Yang. Telling left from right: Identifying geometry-aware semantic correspondence. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3076–3085, 2024. 4
- [52] Jingbo Zhang, Xiaoyu Li, Ziyu Wan, Can Wang, and Jing Liao. Text2nerf: Text-driven 3d scene generation with neural radiance fields. *IEEE Transactions on Visualization and Computer Graphics*, 2024. 3
- [53] Junsheng Zhou, Yu-Shen Liu, and Zhizhong Han. Zero-shot scene reconstruction from single images with deep prior assembly. *arXiv preprint arXiv:2410.15971*, 2024. 3
- [54] Mengqi Zhou, Jun Hou, Chuanchen Luo, Yuxi Wang, Zhaoxiang Zhang, and Junran Peng. Scenex: Procedural controllable large-scale scene generation via large-language models. *arXiv preprint arXiv:2403.15698*, 2024. 3
- [55] Xiaoyu Zhou, Xingjian Ran, Yajiao Xiong, Jinlin He, Zhiwei Lin, Yongtao Wang, Deqing Sun, and Ming-Hsuan Yang. Gala3d: Towards text-to-3d complex scene generation via layout-guided generative gaussian splatting. *arXiv preprint arXiv:2402.07207*, 2024. 3
- [56] Yang Zhou, Zachary White, and Evangelos Kalogerakis. Scenegraphnet: Neural message passing for 3d indoor scene augmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7384–7392, 2019. 3