# AnyHome: Open-Vocabulary Generation of Structured and Textured 3D Homes

Rao Fu[*†1], Zehao Wen[*2], Zichen Liu[*2], and Srinath Sridhar[1]

[1] Brown University
[2] Shenzhen College of International Education
[*]Equal Contribution. [†] Corresponding Author.
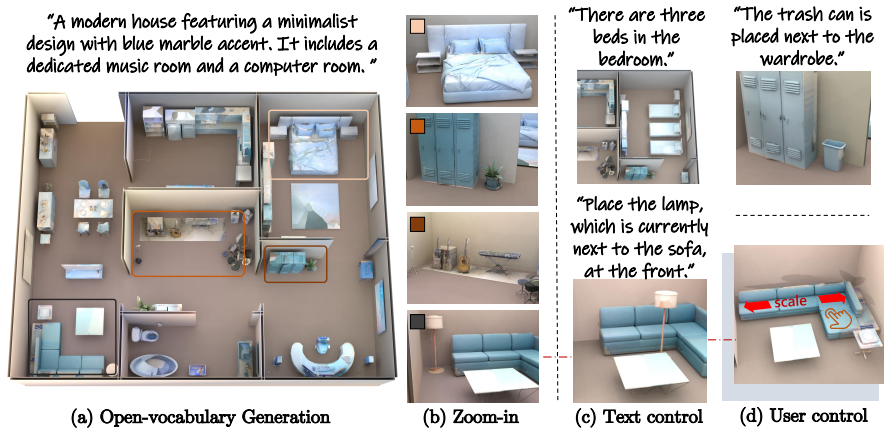`https://ivl.cs.brown.edu/research/anyhome`

**Abstract.** Inspired by cognitive theories, we introduce **AnyHome**, a framework that translates any text into well-structured and textured indoor scenes at a house-scale. By prompting Large Language Models (LLMs) with designed templates, our approach converts provided textual narratives into amodal structured representations. These representations guarantee consistent and realistic spatial layouts by directing the synthesis of a geometry mesh within defined constraints. A Score Distillation Sampling process is then employed to refine the geometry, followed by an egocentric inpainting process that adds lifelike textures to it. **AnyHome** stands out with its editability, customizability, diversity, and realism. The structured representations for scenes allow for extensive editing at varying levels of granularity. Capable of interpreting texts ranging from simple labels to detailed narratives, **AnyHome** generates detailed geometries and textures that outperform existing methods in both quantitative and qualitative measures.

**Keywords:** Scene Generation · 3D Synthesis · Vision and Language

## 1 Introduction

Homes are pivotal to our existence. They are the silent witnesses to our routines and landmark moments, repositories of personal and collective memories that influence our well-being and behavior. Imagine the possibilities if we could articulate our ideal living spaces in natural language and see them come to life. **AnyHome** embodies this vision, offering a framework that transforms free-form, open-vocabulary textual narratives into diverse, house-scale 3D indoor scenes with realistic appearances. These scenes can be used in a variety of domains, including interior design, game development, augmented and virtual reality, as well as the training for embodied agents. They feature intricate structure layouts and naturalistic textures that are readily modifiable, bridging a crucial gap in contemporary digital design practices.

Previous research in text-guided 3D scene generation [17,24,26,65] has made notable advances but often struggles with creating robust 3D structures, sometimes resulting in rooms with open ends or repetitive layouts. Studies focusing on

"A modern house featuring a minimalist design with blue marble accent. It includes a dedicated music room and a computer room."

"There are three beds in the bedroom."

"The trash can is placed next to the wardrobe."

"Place the lamp, which is currently next to the sofa, at the front."

scale

(a) Open-vocabulary Generation    (b) Zoom-in    (c) Text control    (d) User control

**Fig. 1: Example house-scale indoor scene generated by AnyHome.** Users can input any textual description of an indoor scene, and the system is capable of generating house floorplans, room layouts, object placements, and stylistic appearances accordingly. The generated indoor scene is represented by structured and textured room and object meshes. **AnyHome** enables the synthesis of diverse indoor scenes, allowing users to control scene generation at any stage—from textual input and intermediate representation to the generated meshes.

structured scene generation [44,45,48,64,72] are typically confined to predefined room and furniture types, lacking the flexibility to accommodate customization for various room types, furniture pieces, small objects, and their respective arrangements. At the same time, the generation of house-scale scenes with diverse structures is crucial, especially for applications requiring seamless navigation across different rooms, like video games or virtual training for embodied agents. Some studies have explored house-scale scene synthesis [13], but they often rely on predefined artificial textures, which may limit controllability and realism, hindering certain downstream tasks [30].

In response to these limitations, **AnyHome** focuses on creating indoor scenes that are customizable through **open-vocabulary** text inputs, featuring **structured representations**, **realistic texture**, and are **scalable to house-size**. We draw inspiration from two key hypotheses in environmental cognition. The *amodal spatial image* hypothesis suggests that environment recognition transcends sensory modalities like vision or hearing, proposing that individuals apprehend their surroundings through an abstract, symbolic, map-like representation [23,35]. Conversely, the *visual recording* hypothesis posits that visual experiences function similarly to a video camera, perceiving the environment through continuous, egocentric, path-oriented exploration [22,51]. Learning from both concepts, **AnyHome** maintains an amodal, hierarchical geometry representation during the generation process to conceptualize the environment's structure.

To enhance visual realism, we adopt an egocentric exploration approach for inpainting, encouraging the model to detail the environment as it "sees".

Specifically, the scene generation process in **AnyHome** unfolds in several stages: textual input modulation, hierarchical structured geometry generation, along with egocentric refinement and inpainting. This approach is markedly text-controllable. Users can design a scene through free-form text descriptions, which are subsequently converted into modular descriptions for aspects such as floorplans, room layouts, object retrieval and placement, as well as scene appearances, all facilitated by querying Large Langueg Models (LLMs) with specialized templates. These descriptions transformed to graph-based intermediate representations, which fosters customizability and solves the problem of insensible geometry generated directly from the LLMs. Following this, the scene's floorplan is obtained through a generative model [44, 45], and room layouts together with object placements are determined according to LLM-dictated placement rules. We use databases for furniture [19] and objects [12] to populate these scenes. During the egocentric inpainting phase, a camera trajectory imitating first-person exploration is generated. Textures are painted along this trajectory using a depth-aware, text-conditioned inpainting model, which aligns the texture with existing geometry [65]. However, we note inaccuracies in furniture and object placements due to the inherent limitations in object canonicalization from the mesh datasets. To counter, we integrate a Score Distillation Sampling (SDS) process with a differentiable renderer to refine object placement, inspired by recent 3D generation advancements using SDS loss [52] for geometry optimization.

Our experiments demonstrate that **AnyHome** effectively generates house-scale scenes with compelling structures, visually appealing textures, and strong alignment with provided textual inputs. The language accessibility and structured representations enable users to control scene generation at various levels, from room types and layouts to object placements and appearances.

To summarize, our key contributions include:

1. Developing a systematic and reliable approach for creating diverse, text-controlled, texture-realistic, and modifiable scenes at a house-scale, catering to a wide range of applications.
2. Utilizing LLMs to convert open-vocabulary textual input into structured representations, allowing detailed control using language-based customization while maintaining structure consistency.
3. Enhancing object placement using an SDS process, thus increasing the system's robustness and versatility in creating sensible scenes.
4. Innovating an egocentric inpainting process that follows a camera trajectory generated automatically to explore each object in the room.
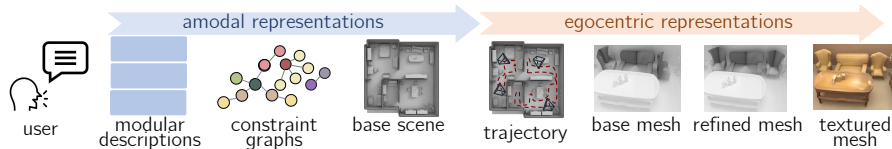
## 2   Related Work

**House-scale Floorplan Generation.** House floorplan generation has been well-studied through methods like shape grammars and iterative generation processes [37, 39, 43, 50]. With the advent of deep learning, the focus has shifted

towards using graph-constrained generative networks [4, 25, 36, 44, 45, 60, 62, 63], accommodating a variety of user-defined constraints. As our work concentrates on the functionality of rooms, we utilize bubble diagrams as a constraint [45], which outline pre-defined room types and the functional connections between them. Our primary contribution lies in synthesizing bubble diagrams based on user textual input and expanding these diagrams to encompass any room type, beyond the predefined categories.

**Room-scale Scene Generation.** Room-scale scene generation involves creating 3D content, such as furniture and smaller objects, and determining their placement. Traditional methods typically start by constructing a set of 3D objects and then optimizing their placement using iterative methods [13, 16, 20], non-linear optimization [15, 54, 75, 79], or manual interaction [26, 40]. However, with the development of large-scale indoor scene datasets [18, 49, 78], recent approaches have shifted towards using generative networks, employing techniques like feed-forward networks [81], VAEs [53], GANs [76], autoregressive models [32, 66–68], and diffusion models [64, 72, 80]. These advanced methods can generate diverse and realistic 3D scenes but often struggle to place objects unpresented in the training dataset. Recent studies [14, 77] have attempted to overcome these limitations by using Large Language Models (LLMs) to generate style sheet languages for object placement. Our method generates room layouts in two steps. We initially create 3D layouts using LLMs in a few-shot manner and then employ a differentiable renderer to refine the placement of objects.

**Text-to-Shape Generation.** Text-to-Shape Generation has seen considerable development, with many studies utilizing feed-forward methods [1, 8, 11, 21, 29, 34, 41, 46, 58, 59, 71] to train generators on 3D data. These methods generate 3D shapes efficiently but their ability to generalize is often limited by the size of the 3D data on which they are trained. Some recent approaches [28, 33, 42, 52, 69] leverage pre-trained visual-language knowledge to optimize 3D representations, showing promising generalization to open-vocabulary textual inputs. Rather than generating shapes from scratch, some research [5, 7, 27, 38, 57] focuses on texture generation, aiming to create textures represented as UV maps from a given mesh. However, optimizing camera views for a single shape is relatively straightforward and cannot directly translate to the complexity of scene geometry. Our paper introduces a novel method for generating camera trajectories specifically for texturing more complex and varied scene structures.

**Text-to-Scene Generation.** Text-to-Scene Generation has followed two primary approaches. The first approach [6, 64, 68] directly uses 3D representations for generation. These methods begin by transforming user text inputs into graph-based representations or lists of shape codes, subsequently generating 3D scenes based on these intermediate representations. While effective in creating structured 3D representations, they often fail to generate realistic textures and unseen objects. The second approach addresses these issues through the image domain. Some studies [10] use panorama images for scene representation and develop end-to-end models that convert text to panorama for scene generation. Other research [2, 3, 17, 24, 61] explores scenes from first-person perspectives, employing

**Fig. 2: Two-stage Generation Process. AnyHome** unfolds two primary steps: First, amodal representations are generated from user's text input, which involves constructing modular text descriptions, constrained graphs, and hierarchical structured base mesh. Following this, the method embarks on an egocentric exploration stage, where a navigation trajectory is generated, enabling the refinement and texturing of the base mesh from different viewpoints.
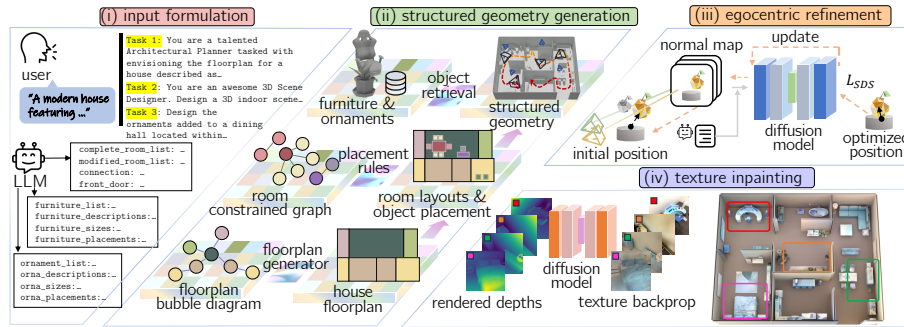
large-scale text-conditioned image generation models for appearance creation. Techniques like neural density fields or depth are used to ensure 3D consistency. This approach excels in generating diverse appearances but struggles with creating structured 3D representations and consistent shapes. Our method overcomes these challenges by adopting a two-stage generation process. We first generate spatial layouts in a few-shot manner and then inpaint textures using first-person viewpoints, which guarantees both structured outputs and diverse appearances.

## 3   AnyHome

Fig. 2 illustrates the two-stage generation process of **AnyHome**. Our approach incorporates insights from environmental cognition hypothesis. We adopt the *amodal spatial image* hypothesis [23, 35] to represent houses with an amodal hierachical structured representation. Drawing on the *visual recording* hypothesis [22, 51], we further refine and enrich the visual details of houses using egocentric inpainting. Fig. 3 illustrates the main components of the pipeline. Our framework comprises three main components: First, we modulate the textual input using Large Language Models (LLMs) as detailed in Sec. 3.1, which involves comprehending and elaborating user input. Second, in Sec. 3.2, we focus on synthesizing graph-based structured representations that are important for controllability, and use them to generate the base geometry. Finally, Sec. 3.3 describes our process of refinement and inpainting through egocentric exploration, which refines object placements and adds texture, enhancing the realism of the generated scene.

### 3.1   Textual Input Modulation.

This initial component of our framework performs two key functions. First, it allows users to offer an unrestricted description of the scene to guide and control the generation of designs. Second, it utilizes the common-sense knowledge inherent within LLMs [70, 82] to convert the provided input into modular descriptions

**Fig. 3: Pipeline.** Taking a free-form textual input, our pipeline generates the house-scale scene by: (i) comprehending and elaborating on the user's textual input through querying an LLM with templated prompts; (ii) converting textual descriptions into base geometry using structured intermediate representations; (iii) employing an SDS process with a differentiable renderer to refine object placements; and (iv) applying depth-conditioned texture inpainting for egocentric texture generation.

for house floorplans, room layouts, and appearances of furnishings and ornaments. Notably, LLMs are instructed to render graph-like structures that act as intermediate constraints for the generation of house floorplans and room layouts. We have designed three sets of prompts to guide the LLMs: for house floorplan generation ($p_{\text{floorplan}}, p_{\text{map}}$), for room layout and object placement ($p_{\text{room}}$), and for room appearance ($p_{\text{appearance}}$). Fig. 3 (i) provides an example of using ($p_{\text{room}}$) to generate room layout descriptions. The design of these prompts bases on three principles: format compatibility with subsequent modules, adherence to user's textual specifications, and maximal detail elaboration by the LLM-agent, which is encouraged to adopt the perspective of an interior designer or navigator. Additionally, numerical data provided to LLMs is converted to a real-world scale (in meters) to fully utilize their knowledge base. Further details on these prompts can be found in the Supplementary Material.

### 3.2   Hierarchical Structured Geometry Generation.

In this phase, we concentrate on converting modular descriptions into a detailed base geometry of the scene, covering elements like floors, walls, furniture, and smaller objects. For precise spatial control, we employ two graph-based intermediate representations: one dedicated to the floorplan and the other to room layouts. Fig. 3 (ii) illustrates this hierarchical generation process.

**House Floorplan Generation.** A floorplan is crucial as it outlines the spatial arrangement, dimensions, and functionalities of rooms. While zero-shot methods like LayoutGPT [14] can generate a variety of room types, they sometimes yield impractical configurations, as illustrated in Fig. 8. To address these challenges,

we initially create bubble diagrams, a graph-based representation similar to that in *HouseGAN++* [45], and input them into pre-trained floorplan generation networks to facilitate text-controlled floorplan synthesis.

In this approach, each floorplan is represented by a bubble-diagram $G(N, E)$, where nodes $N$ symbolize room types and edges $E$ represent functional connections. We query an LLM with a specific floorplan generation prompt $p_{\text{floorplan}}$ to convert user inputs into these diagrams. Utilizing the pre-trained, graph-conditioned floorplan generator *HouseGAN++*, we transform these diagrams into a set of masks $M$ that detail the spatial layout.

Considering *HouseGAN++*'s limitations with its *RPLAN* training data, we introduce additional prompts $p_{\text{map}}$ to map unconventional room types to *RPLAN*'s recognized categories. Upon defining $M$, we construct the house's base mesh at a defined height $h$. This method, combining LLMs with a specialized floorplan network, maintains variety while ensuring spatial realism.

**Room Layout and Object Placement Generation.** With the house floorplan established, we advance to generating room layouts and object placements. While existing methods have shown limitations with unique room types [48, 72] and simplistic layouts [14] as evidenced in Fig. 8, our approach introduces a constrained layout graph $\bar{G}_i(\bar{N}_i, \bar{E}_i)$ to enhance coherence and flexibility. This graph captures dimensions and descriptions of objects in its nodes $\bar{N}_i$ and the relationships between objects in edges $\bar{E}_i$.

Using LLMs, we generate a graph $\bar{G}_i$ for each room based on its area and type, as well as the user input. LLMs struggle to envision a room with a diverse array of objects and to manage their complex interrelations simultaneously, and hence we prompt the LLM twice during this graph generation process. We employ two distinct prompt templates: $p_{\text{furniture}}$ for large furniture and $p_{\text{ornament}}$ for smaller ornaments, to address these elements separately.

The features stored in the nodes and edges are generated in a step-by-step manner. For instance, in furniture placement, we initially generate a list of potential furniture items along with their dimensions. Subsequently, drawing on the concept of Semantic Asset Group (SAG) [13], the LLM is prompted to identify the sub-graphs $C_i$ within $\bar{G}_i$ that represent groups of commonly associated items and designate an anchor piece within each group, such as a table in a set of a table and chairs.

With a predefined set of placement rules like *place_corner(·)*, *place_beside(·)*, and *place_wall(·)*, we direct the LLM to determine the placement rules between the anchor object and other items within the SAG. The corresponding placement algorithm for each rule then dictates the objects' locations, orientations, and adjustments in the event of obstructions, thus optimizing space usage. The Supplementary Material provides a detailed account of all placement rules, including their functions and parameters. This stage results in a series of bounding boxes $B_i$ outlining the final layout.

**Object Retrieval.** To enhance the scene with appropriate furniture and decor, we utilize comprehensive mesh datasets [12,19], retrieving items that best match the LLM-generated descriptions. By employing CLIP [55] and Sentence Transformers [56], we ensure that the retrieved items align closely with the textual descriptions, seamlessly integrating them into the overall house mesh.

### 3.3    Egocentric Refinement and Inpainting.

In this stage, we enhance the base geometry in an egocentric manner, refining its details and layouts as well as adding textures. This process is illustrated in Fig. 3 (iii) and (iv). We start by creating an egocentric trajectory, followed by layout refinement using Score Distillation Sampling (SDS) loss. The final step involves using a depth-conditioned inpainting model for texturing, ensuring coherence between the geometry and textures. Also, a differentiable renderer is employed to maintain texture consistency from multiple viewpoints.

**Generating Egocentric Trajectories.** Our trajectory generation adheres to three principles: (1) ensuring complete coverage by guiding the camera along walls; (2) focusing on object-centric views, facilitated by our structured geometry; and (3) implementing double traversal to capture views both towards and away from the interior walls. Additionally, we include random camera samples to achieve comprehensive coverage. This process begins by creating an obstacle field based on the base geometry. By setting a specific threshold, we can identify key points within the field at certain magnitudes to define the trajectory, with the gradient indicating the camera's direction. For random sampling, we select from a set of closed-loop trajectories as feature points and employ a Bezier curve to smoothly interpolate the camera path. Detailed information on the algorithm and its visualizations can be found in the Supplementary Material.

**Refinement and Inpainting with Text-to-Image Models.** With the egocentric camera trajectories set, we still face challenges: un-textured geometry and layout inaccuracies resulting from placement rules and inconsistent object coordinate systems from the mesh datasets. We address these by optimizing object placement and texture with image domain losses. In a given scene with $N$ meshes, denoted as $M = \{m_i | i \in [1, N]\}$, the position of each mesh $P = \{p_i | i \in [1, N]\}$ is defined relative to its anchor mesh, parameterized as the relative translation and quaternion rotation. The corresponding vertex colors for these meshes are symbolized by $M^c = \{m_i^c | i \in [1, N]\}$. Following Fantasia3D [9], we divide the process into two distinct stages: placement refinement and mesh texturing.

   **Structure Refinement.** We render the mesh set's normal map $(n, o)$ and mask $(o)$ from camera views $(c)$ using a differentiable renderer $(\Psi)$. The position parameters $P$ are then refined using the SDS loss, calculated as follows:

$$\nabla_P \mathcal{L}_{\text{SDS}}(\phi, \Psi(M, c)) = \mathbb{E}_{t,\epsilon}\left[w(t)(\hat{\epsilon}_\phi(\tilde{n}_t; y, t) - \epsilon)\frac{\partial z}{\partial P}\right], \qquad (1)$$

**Fig. 4: Open-Vocabulary Generation Results.** Top: Input text prompt. Middle: Bird's-eye view of the scenes. Bottom: Egocentric view of the scenes. **AnyHome** interprets users' textual inputs and produces structured scenes with realistic textures. It can create a serene and culturally rich environment (Left - "Japanese tea house"), render a more dramatic and stylized ambiance (Middle - "haunted house"), and synthesize unique house types (Right - "cat cafe").

where $\phi$ parameterizes the pre-trained stable diffusion model, and $\hat{\epsilon}_\phi(z_t^{\tilde{n}}; y, t)$ represents the noisy function given the noisy image $\tilde{n}_t$, text embedding $y$, and noise level $t$, with $\epsilon$ being the added noise. We utilize mutli-view images with a batch size of 8 for optimization. This refinement improves object positioning, surface adherence, as well as penetration issues.

**Texturing with Depth-conditioned Inpainting.** To ensure realistic generation and preservation of 3D geometry, depth maps ($d_i$) are rendered from each viewpoint. The LLM transforms user textual input ($t$) into an appearance diffusion prompt ($p_{\text{diffusion}}$) to the diffusion model, utilizing a predefined appearance prompt ($p_{\text{appearance}}$). This diffusion prompt directs a depth-to-image inpainting model to produce textured images ($I_i = \phi(p_{\text{diffusion}}, d_i)$). These images are then back-projected into 3D space, iteratively refining the mesh texture as $M^{c_{i+1}} = \Psi(M^{c_i}, I_i)$. This approach ensures the alignment of the textured mesh with the original 3D geometry and its detailed design elements, such as adding calligraphy to scrolls.

## 4    Results

In this section, we present our results on open-vocabulary 3D house generation and editing. We also provide comparison with other methods both quantitatively and qualitatively.

**Fig. 5: Open-vocabulary Editing Results.** Examples showcase **AnyHome**'s capability to modify room types, layouts, object appearances, and overall design through free-form user input. **AnyHome** also supports comprehensive style alterations and sequential edits, all made possible by its hierarchical structured geometric representation and robust text controllability.
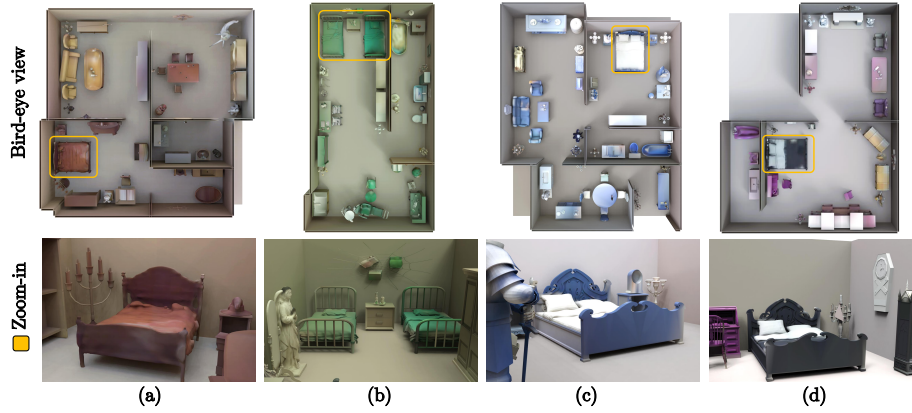
### 4.1   Open-vocabulary Scene Generation.

Fig. 4 illustrates **AnyHome**'s ability to generate house-scale 3D scenes from open-vocabulary, free-form text. The framework interprets and elaborates on a wide array of styles, including both eastern and western designs. This allows for the creation of stylistically coherent rooms with appropriately placed objects. One notable aspect of **AnyHome** is its capacity to generate an extended range of room types, such as tea rooms, moving beyond the limitations of the standard RPLAN list [74]. The showcased examples are a testament to **AnyHome**'s proficiency in creating various stylistic scenes, each distinguished by their coherent floorplans, layouts, and the high-quality, consistent textures that resonate with the specified textual descriptions.

### 4.2   Open-vocabulary Editing.

We present a series of detailed editing results to further illustrate the customizability of **AnyHome** in Fig. 5. These examples highlight the framework's text-controllability across various levels: from altering room types (left row 1 - from kitchen to study room), adjusting room layouts (left row 2 - from cushion set to a Buddha statue), modifying room styles (middle row 1 - from haunted to Barbie pink), to changing object appearances (right row 2 - from green seats to orange seats). Additionally, **AnyHome** supports free-form design modifications (right row 1 -from reading to movie area) and facilitates sequential editing (middle row

2). These functionalities are made possible by the system's modular text prompt design and its hierarchical, structured geometry representations.



**Fig. 6: Diverse Scene Results.** Four distinct scenes generated for the prompt `"A one-bedroom, one-bathroom haunted house featuring dark wood and antique furnishings."` **AnyHome** produces houses with diverse floorplans, room types, room layouts, objects and textures.
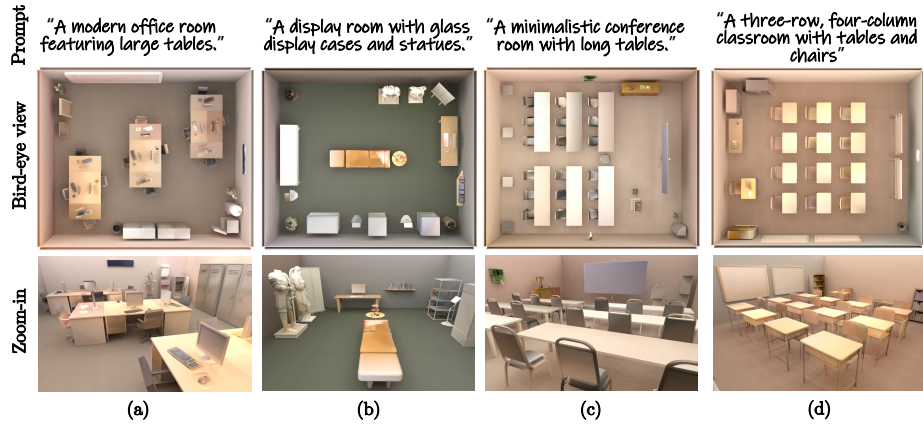
### 4.3 Diversity.

Fig. 6 showcases four unique scenes generated by **AnyHome** from a single text input, demonstrating its capability for diversity. **AnyHome** successfully generates varied floorplans (a five-room scene in **a** and a three-room scene in **c**), room types (a laboratory in **b**), room layouts (two parallel beds in **b**), room objects (a soldier in **c**), and appearances (different styles and colors), all while maintaining coherence with the same textual description. Fig. 7 highlights **AnyHome**'s ability to create layouts beyond traditional "home-like" environments. The system can produce various room types such as offices, display rooms, conference rooms, and classrooms with neatly symmetrical patterns.

### 4.4 Comparison.

We feature quantitative and qualitative comparisons of **AnyHome** against other baseline methods in the following section. We provide comparisons with concurrent works in the Supplementary Material.

**Quantitative Comparison with Baselines** To evaluate the effectiveness of **AnyHome** in scene layout generation and text-to-scene alignment, we conduct a comparative analysis using several key metrics. **(1)** *Out of Bound Rates (OOB):*
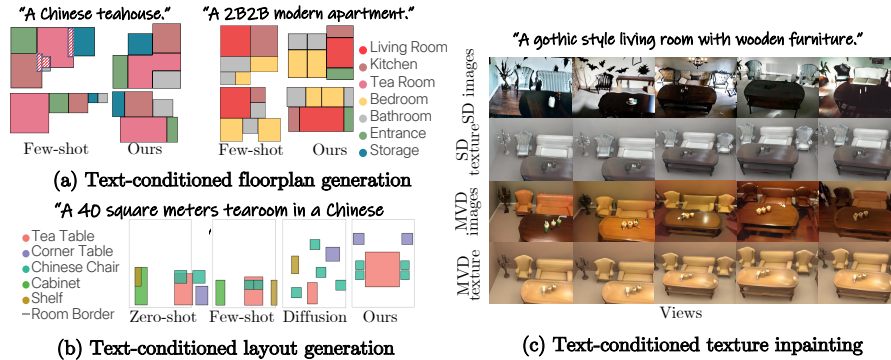
Fig. 7: **Diverse Layout Results. AnyHome** can generate diverse layouts beyond usual room types, showcasing its versatility to be used for different applications. The layouts synthesized can be neatly symmetrical and orderly, echoing to the specifics provided in the prompt.

The quality of layouts is quantified by measuring the frequency of objects intersecting with each other and the walls, or extending beyond room boundaries, following the approach in LayoutGPT [14]. **(2)** *Caption Similarity (Caption-sim):* This metric assesses the alignment between text and scene content. It involves generating captions [31] for rendered frames and computing sentence similarity [56] with the diffusion model $p_{\text{diffusion}}$. This is crucial for evaluating how well the layout and appearance match the user's input from an egocentric viewpoint. **(3)** *CLIP Similarity (CLIP-sim):* The correlation between text and scene content is also assessed using CLIP Similarity. This involves extracting image features via CLIP from rendered frames and calculating the CLIP-Score relative to the text features of the input using the CLIP/H-14 model. **(4)** *CLIP Style Similarity (CLIP-style-sim):* Similar to CLIP Similarity, this metric focuses on evaluating the correspondence between the scene's overall style and the user's original text input.

We compare **AnyHome** against several baselines. **LayoutGPT+Retrieval** utilizes LayoutGPT [14], a method that employs a mix of few-shot and zero-shot methods for object placement, with objects subsequently retrieved from the 3D-Front [18] furniture and Objaverse [12] databases, retaining their original texture. **CG+Retrieval** employs our constrained graph-based furniture placement algorithm for layout, followed by object retrieval. **CG+Inpainting** combines our layout approach with MVDiffusion [65] for texturing. **AnyHome** enriches these with SDS-based layout refinement.

Tab. 1 presents the quantitative evaluation. Comparing the first and second rows, using our constrained graph-based representation (CG) significantly improves layout generation over the few-shot LLM-based method without in-

**Fig. 8: Comparison of Structure Generation and Texture Inpainting.** For structure generation, our method generates more complex and reasonable floorplans and layouts. For texture inpainting, our method generates view-consistent and text-content aligned texture.

**Table 1:** Comparison of layout and content generation quality between **AnyHome** and other baselines using Out-of-Boundary Rate (OBB), Caption Similarity (Caption-*sim*), CLIP Similarity (CLIP-*sim*), and CLIP Style Similarity (CLIP-style-*sim*). Mean and variance from 10 runs with random viewpoints are reported. **AnyHome** outperforms all baseline methods.

| Method | OOB | Caption-*sim* | CLIP-*sim* | CLIP-style-*sim* |
|---|---|---|---|---|
| LayoutGPT+Retrieval | 69.4 | 9.7 | 8.3 | 6.4 |
| CG+Retrieval | 34.2 | 11.2 | 11.5 | 7.0 |
| CG+Inpainting | 34.2 | 15.9 | 27.9 | **17.5** |
| **AnyHome** | **23.7** | **16.2±0.1** | **29.3±0.6** | **17.5±0.5** |

termediate representations (LayoutGPT), as evidenced by a lower OOB rate and higher content generation scores. Comparing the second and third rows, the egocentric inpainting process further improves text-content and text-style correspondences, elevating these scores (Caption-*sim*, CLIP-*sim*, CLIP-style-*sim*). Comparing the third and fourth rows, the refinement stage further enhances the quality of the layout by reducing overlaps and refining object placements, as indicated by improved scores across all metrics. We include random viewpoint perturbations in the evaluation, such as z-axis rotation and ground translation. The low variances demonstrate the minimal impact of viewpoints on the results.

**Qualitative Comparison for Structure Generation.** In Fig. 8 panels (a) and (b), we compare our method with existing techniques for floorplan and room layout generation, including zero-shot LLM [47], few-shot LLM [14], and diffusion-based methods [72]. For house floorplan generation, LLM-based methods often produce rooms with intersecting walls, overly simplistic structures, or illogical configurations, such as a bedroom situated within a bathroom. Our method surpasses these direct LLM-generated plans, especially with abstract

prompts, by preserving room relationships and accommodating diverse shapes and sizes. This approach effectively addresses open-vocabulary input scenarios. For room layout generation, our SAG-based method demonstrates a clear advantage over traditional few-shot LLMs and diffusion methods. It results in more logically arranged furniture, avoiding common issues such as overlapping bounding boxes and spatial incoherence, particularly in open-vocabulary settings.

**Qualitative Comparison for Texture Inpainting.**   Fig. 8 (c) presents a qualitative analysis of our texture inpainting approach in comparison with other existing methods, guided by the text prompt `"A gothic living room with wooden furniture."` In this assessment, we evaluate baseline methods such as Stable-Diffusion-2 with depth conditioning (SD images) and the results of back-propagating image pixel colors to mesh vertices using differentiable rendering (SD texture). Additionally, we examine images generated by MVDiffusion (MVD images) and their subsequent mesh vertex back-propagation (MVD texture). While SD images exhibit richly detailed textures, they often lead to inconsistencies across different views, which undermines the effectiveness of back-propagation. In contrast, MVD images demonstrate a higher level of view consistency. The back-propagation of MVD images (MVD texture) further enhances this consistency across various viewpoints. By integrating this method, which leverages the strengths of MVDiffusion and differentiable rendering, we achieve textures that are not only consistent but also smoothly transitioned, accurately reflecting the textual prompt provided.

## 5   Conclusion

In this paper, we present **AnyHome**, a novel framework designed to generate house-scale scenes from open-vocabulary textual inputs. Leveraging specific prompt templates, our system effectively interprets, follows, and enriches user-provided text to create detailed scenes. Utilizing a two-stage generation approach, which combines LLM-designated rules with SDS loss refinement, our method proficiently arranges a diverse array of objects in a realistic manner. The resulting houses feature a structured geometry representation, enhancing the ease of user editing and modification. We believe that **AnyHome** paves the way for a wide range of applications, including complex 3D interior design, immersive augmented and virtual reality experiences, dynamic gaming environments, and advanced training modules for embodied agents.

**Limitations and Future Work.** Still, **AnyHome** faces limitations due to the current LLMs' understanding of 3D spaces and the challenges of maintaining consistency in multi-view inpainting using existing techniques. Future work will focus on overcoming these challenges to improve the system's capacity for generating detailed and coherent 3D environments from textual inputs.

## Acknowledgements

## References

1. Achlioptas, P., Diamanti, O., Mitliagkas, I., Guibas, L.: Learning representations and generative models for 3d point clouds. In: International conference on machine learning. pp. 40–49. PMLR (2018)
2. Bahmani, S., Park, J.J., Paschalidou, D., Yan, X., Wetzstein, G., Guibas, L., Tagliasacchi, A.: Cc3d: Layout-conditioned generation of compositional 3d scenes. arXiv preprint arXiv:2303.12074 (2023)
3. Bautista, M.A., Guo, P., Abnar, S., Talbott, W., Toshev, A., Chen, Z., Dinh, L., Zhai, S., Goh, H., Ulbricht, D., et al.: Gaudi: A neural architect for immersive 3d scene generation. Advances in Neural Information Processing Systems **35**, 25102–25116 (2022)
4. Bisht, S., Shekhawat, K., Upasani, N., Jain, R.N., Tiwaskar, R.J., Hebbar, C.: Transforming an adjacency graph into dimensioned floorplan layouts. In: Computer Graphics Forum. vol. 41, pp. 5–22. Wiley Online Library (2022)
5. Cao, T., Kreis, K., Fidler, S., Sharp, N., Yin, K.: Texfusion: Synthesizing 3d textures with text-guided image diffusion models. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 4169–4181 (2023)
6. Chang, A.X., Eric, M., Savva, M., Manning, C.D.: Sceneseer: 3d scene design with natural language. arXiv preprint arXiv:1703.00050 (2017)
7. Chen, D.Z., Siddiqui, Y., Lee, H.Y., Tulyakov, S., Nießner, M.: Text2tex: Text-driven texture synthesis via diffusion models. arXiv preprint arXiv:2303.11396 (2023)
8. Chen, K., Choy, C.B., Savva, M., Chang, A.X., Funkhouser, T., Savarese, S.: Text2shape: Generating shapes from natural language by learning joint embeddings. In: Computer Vision–ACCV 2018: 14th Asian Conference on Computer Vision, Perth, Australia, December 2–6, 2018, Revised Selected Papers, Part III 14. pp. 100–116. Springer (2019)
9. Chen, R., Chen, Y., Jiao, N., Jia, K.: Fantasia3d: Disentangling geometry and appearance for high-quality text-to-3d content creation. arXiv preprint arXiv:2303.13873 (2023)
10. Chen, Z., Wang, G., Liu, Z.: Text2light: Zero-shot text-driven hdr panorama generation. ACM Transactions on Graphics (TOG) **41**(6), 1–16 (2022)
11. Cheng, Y.C., Lee, H.Y., Tulyakov, S., Schwing, A.G., Gui, L.Y.: Sdfusion: Multi-modal 3d shape completion, reconstruction, and generation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4456–4465 (2023)
12. Deitke, M., Schwenk, D., Salvador, J., Weihs, L., Michel, O., VanderBilt, E., Schmidt, L., Ehsani, K., Kembhavi, A., Farhadi, A.: Objaverse: A universe of annotated 3d objects. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 13142–13153 (2023)
13. Deitke, M., VanderBilt, E., Herrasti, A., Weihs, L., Ehsani, K., Salvador, J., Han, W., Kolve, E., Kembhavi, A., Mottaghi, R.: Procthor: Large-scale embodied ai using procedural generation. Advances in Neural Information Processing Systems **35**, 5982–5994 (2022)

14. Feng, W., Zhu, W., Fu, T.j., Jampani, V., Akula, A., He, X., Basu, S., Wang, X.E., Wang, W.Y.: Layoutgpt: Compositional visual planning and generation with large language models. arXiv preprint arXiv:2305.15393 (2023)
15. Fisher, M., Ritchie, D., Savva, M., Funkhouser, T., Hanrahan, P.: Example-based synthesis of 3d object arrangements. ACM Transactions on Graphics (TOG) **31**(6), 1–11 (2012)
16. Fisher, M., Savva, M., Li, Y., Hanrahan, P., Nießner, M.: Activity-centric scene synthesis for functional 3d scene modeling. ACM Transactions on Graphics (TOG) **34**(6), 1–13 (2015)
17. Fridman, R., Abecasis, A., Kasten, Y., Dekel, T.: Scenescape: Text-driven consistent scene generation. arXiv preprint arXiv:2302.01133 (2023)
18. Fu, H., Cai, B., Gao, L., Zhang, L.X., Wang, J., Li, C., Zeng, Q., Sun, C., Jia, R., Zhao, B., et al.: 3d-front: 3d furnished rooms with layouts and semantics. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 10933–10942 (2021)
19. Fu, H., Jia, R., Gao, L., Gong, M., Zhao, B., Maybank, S., Tao, D.: 3d-future: 3d furniture shape with texture. International Journal of Computer Vision pp. 1–25 (2021)
20. Fu, Q., Chen, X., Wang, X., Wen, S., Zhou, B., Fu, H.: Adaptive synthesis of indoor scenes via activity-associated object relation graphs. ACM Transactions on Graphics (TOG) **36**(6), 1–13 (2017)
21. Fu, R., Zhan, X., Chen, Y., Ritchie, D., Sridhar, S.: Shapecrafter: A recursive text-conditioned 3d shape generation model. Advances in Neural Information Processing Systems **35**, 8882–8895 (2022)
22. Gibson, J.J.: The ecological approach to visual perception: classic edition. Psychology press (2014)
23. Giudice, N.A.: 15. navigating without vision: Principles of blind spatial cognition. Handbook of behavioral and cognitive geography p. 260 (2018)
24. Höllein, L., Cao, A., Owens, A., Johnson, J., Nießner, M.: Text2room: Extracting textured 3d meshes from 2d text-to-image models. arXiv preprint arXiv:2303.11989 (2023)
25. Hu, R., Huang, Z., Tang, Y., Van Kaick, O., Zhang, H., Huang, H.: Graph2plan: Learning floorplan generation from layout graphs. ACM Transactions on Graphics (TOG) **39**(4), 118–1 (2020)
26. Huang, I., Krishna, V., Atekha, O., Guibas, L.: Aladdin: Zero-shot hallucination of stylized 3d assets from abstract scene descriptions. arXiv preprint arXiv:2306.06212 (2023)
27. Hwang, I., Kim, H., Kim, Y.M.: Text2scene: Text-driven indoor scene stylization with part-aware details. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 1890–1899 (2023)
28. Jain, A., Mildenhall, B., Barron, J.T., Abbeel, P., Poole, B.: Zero-shot text-guided object generation with dream fields.(2022) (2022)
29. Jun, H., Nichol, A.: Shap-e: Generating conditional 3d implicit functions. arXiv preprint arXiv:2305.02463 (2023)
30. Khanna, M., Mao, Y., Jiang, H., Haresh, S., Schacklett, B., Batra, D., Clegg, A., Undersander, E., Chang, A.X., Savva, M.: Habitat synthetic scenes dataset (hssd-200): An analysis of 3d scene scale and realism tradeoffs for objectgoal navigation. arXiv preprint arXiv:2306.11290 (2023)
31. Li, J., Li, D., Savarese, S., Hoi, S.: Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. arXiv preprint arXiv:2301.12597 (2023)

32. Li, M., Patil, A.G., Xu, K., Chaudhuri, S., Khan, O., Shamir, A., Tu, C., Chen, B., Cohen-Or, D., Zhang, H.: Grains: Generative recursive autoencoders for indoor scenes. ACM Transactions on Graphics (TOG) **38**(2), 1–16 (2019)
33. Lin, C.H., Gao, J., Tang, L., Takikawa, T., Zeng, X., Huang, X., Kreis, K., Fidler, S., Liu, M.Y., Lin, T.Y.: Magic3d: High-resolution text-to-3d content creation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 300–309 (2023)
34. Liu, Z., Wang, Y., Qi, X., Fu, C.W.: Towards implicit text-guided 3d shape generation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 17896–17906 (2022)
35. Loomis, J.M., Lippa, Y., Klatzky, R.L., Golledge, R.G.: Spatial updating of locations specified by 3-d sound and spatial language. Journal of Experimental Psychology: Learning, Memory, and Cognition **28**(2), 335 (2002)
36. Luo, Z., Huang, W.: Floorplangan: Vector residential floorplan adversarial generation. Automation in Construction **142**, 104470 (2022)
37. Ma, C., Vining, N., Lefebvre, S., Sheffer, A.: Game level layout from design specification. In: Computer Graphics Forum. vol. 33, pp. 95–104. Wiley Online Library (2014)
38. Ma, Y., Zhang, X., Sun, X., Ji, J., Wang, H., Jiang, G., Zhuang, W., Ji, R.: X-mesh: Towards fast and accurate text-driven 3d stylization via dynamic textual guidance. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 2749–2760 (2023)
39. Merrell, P., Schkufza, E., Koltun, V.: Computer-generated residential building layouts. In: ACM SIGGRAPH Asia 2010 papers, pp. 1–12 (2010)
40. Merrell, P., Schkufza, E., Li, Z., Agrawala, M., Koltun, V.: Interactive furniture layout using interior design guidelines. ACM transactions on graphics (TOG) **30**(4), 1–10 (2011)
41. Mittal, P., Cheng, Y.C., Singh, M., Tulsiani, S.: Autosdf: Shape priors for 3d completion, reconstruction and generation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 306–315 (2022)
42. Mohammad Khalid, N., Xie, T., Belilovsky, E., Popa, T.: Clip-mesh: Generating textured meshes from text using pretrained image-text models. In: SIGGRAPH Asia 2022 conference papers. pp. 1–8 (2022)
43. Müller, P., Wonka, P., Haegler, S., Ulmer, A., Van Gool, L.: Procedural modeling of buildings. In: ACM SIGGRAPH 2006 Papers, pp. 614–623 (2006)
44. Nauata, N., Chang, K.H., Cheng, C.Y., Mori, G., Furukawa, Y.: House-gan: Relational generative adversarial networks for graph-constrained house layout generation. In: Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16. pp. 162–177. Springer (2020)
45. Nauata, N., Hosseini, S., Chang, K.H., Chu, H., Cheng, C.Y., Furukawa, Y.: House-gan++: Generative adversarial layout refinement network towards intelligent computational agent for professional architects. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 13632–13641 (2021)
46. Nichol, A., Jun, H., Dhariwal, P., Mishkin, P., Chen, M.: Point-e: A system for generating 3d point clouds from complex prompts. arXiv preprint arXiv:2212.08751 (2022)
47. OpenAI: Gpt-4 technical report (2023)
48. Paschalidou, D., Kar, A., Shugrina, M., Kreis, K., Geiger, A., Fidler, S.: Atiss: Autoregressive transformers for indoor scene synthesis. Advances in Neural Information Processing Systems **34**, 12013–12026 (2021)

49. Paschalidou, D., Kar, A., Shugrina, M., Kreis, K., Geiger, A., Fidler, S.: Atiss: Autoregressive transformers for indoor scene synthesis. Advances in Neural Information Processing Systems **34**, 12013–12026 (2021)
50. Peng, C.H., Yang, Y.L., Wonka, P.: Computing layouts with deformable templates. ACM Transactions on Graphics (TOG) **33**(4), 1–11 (2014)
51. Pick, H.L.: Visual coding of nonvisual spatial information. (1974)
52. Poole, B., Jain, A., Barron, J.T., Mildenhall, B.: Dreamfusion: Text-to-3d using 2d diffusion. arXiv preprint arXiv:2209.14988 (2022)
53. Purkait, P., Zach, C., Reid, I.: Sg-vae: Scene grammar variational autoencoder to generate new indoor scenes. In: European Conference on Computer Vision. pp. 155–171. Springer (2020)
54. Qi, S., Zhu, Y., Huang, S., Jiang, C., Zhu, S.C.: Human-centric indoor scene synthesis using stochastic grammar. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5899–5908 (2018)
55. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: International conference on machine learning. pp. 8748–8763. PMLR (2021)
56. Reimers, N., Gurevych, I.: Sentence-bert: Sentence embeddings using siamese bert-networks (11 2019), `http://arxiv.org/abs/1908.10084`
57. Richardson, E., Metzer, G., Alaluf, Y., Giryes, R., Cohen-Or, D.: Texture: Text-guided texturing of 3d shapes. arXiv preprint arXiv:2302.01721 (2023)
58. Sanghi, A., Chu, H., Lambourne, J.G., Wang, Y.a.R.: Clip-forge: Towards zero-shot text-to-shape generation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 18603–18613 (2022)
59. Sanghi, A., Fu, R., Liu, V., Willis, K.D., Shayani, H., Khasahmadi, A.H., Sridhar, S., Ritchie, D.: Clip-sculptor: Zero-shot generation of high-fidelity and diverse shapes from natural language. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 18339–18348 (2023)
60. Shabani, M.A., Hosseini, S., Furukawa, Y.: Housediffusion: Vector floorplan generation via a diffusion model with discrete and continuous denoising. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5466–5475 (2023)
61. Song, L., Cao, L., Xu, H., Kang, K., Tang, F., Yuan, J., Zhao, Y.: Roomdreamer: Text-driven 3d indoor scene synthesis with coherent geometry and texture. arXiv preprint arXiv:2305.11337 (2023)
62. Sun, J., Wu, W., Liu, L., Min, W., Zhang, G., Zheng, L.: Wallplan: synthesizing floorplans by learning to generate wall graphs. ACM Transactions on Graphics (TOG) **41**(4), 1–14 (2022)
63. Tang, H., Zhang, Z., Shi, H., Li, B., Shao, L., Sebe, N., Timofte, R., Van Gool, L.: Graph transformer gans for graph-constrained house generation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 2173–2182 (2023)
64. Tang, J., Nie, Y., Markhasin, L., Dai, A., Thies, J., Nießner, M.: Diffuscene: Scene graph denoising diffusion probabilistic model for generative indoor scene synthesis. arXiv preprint arXiv:2303.14207 (2023)
65. Tang, S., Zhang, F., Chen, J., Wang, P., Furukawa, Y.: Mvdiffusion: Enabling holistic multi-view image generation with correspondence-aware diffusion. arXiv preprint arXiv:2307.01097 (2023)

66. Wang, K., Lin, Y.A., Weissmann, B., Savva, M., Chang, A.X., Ritchie, D.: Planit: Planning and instantiating indoor scenes with relation graph and spatial prior networks. ACM Transactions on Graphics (TOG) **38**(4), 1–15 (2019)
67. Wang, K., Savva, M., Chang, A.X., Ritchie, D.: Deep convolutional priors for indoor scene synthesis. ACM Transactions on Graphics (TOG) **37**(4), 1–14 (2018)
68. Wang, X., Yeshwanth, C., Nießner, M.: Sceneformer: Indoor scene generation with transformers. In: 2021 International Conference on 3D Vision (3DV). pp. 106–115. IEEE (2021)
69. Wang, Z., Lu, C., Wang, Y., Bao, F., Li, C., Su, H., Zhu, J.: Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation. arXiv preprint arXiv:2305.16213 (2023)
70. Wei, J., Tay, Y., Bommasani, R., Raffel, C., Zoph, B., Borgeaud, S., Yogatama, D., Bosma, M., Zhou, D., Metzler, D., Chi, E.H., Hashimoto, T., Vinyals, O., Liang, P., Dean, J., Fedus, W.: Emergent abilities of large language models. Transactions on Machine Learning Research (2022), `https://openreview.net/forum?id=yzkSU5zdwD`, survey Certification
71. Wei, J., Wang, H., Feng, J., Lin, G., Yap, K.H.: Taps3d: Text-guided 3d textured shape generation from pseudo supervision. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 16805–16815 (2023)
72. Wei, Q.A., Ding, S., Park, J.J., Sajnani, R., Poulenard, A., Sridhar, S., Guibas, L.: Lego-net: Learning regular rearrangements of objects in rooms. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 19037–19047 (2023)
73. Wu, T., Yang, G., Li, Z., Zhang, K., Liu, Z., Guibas, L., Lin, D., Wetzstein, G.: Gpt-4v (ision) is a human-aligned evaluator for text-to-3d generation. arXiv preprint arXiv:2401.04092 (2024)
74. Wu, W., Fu, X.M., Tang, R., Wang, Y., Qi, Y.H., Liu, L.: Data-driven interior plan generation for residential buildings. ACM Transactions on Graphics (SIGGRAPH Asia) **38**(6) (2019)
75. Xu, K., Chen, K., Fu, H., Sun, W.L., Hu, S.M.: Sketch2scene: Sketch-based co-retrieval and co-placement of 3d models. ACM Transactions on Graphics (TOG) **32**(4), 1–15 (2013)
76. Yang, M.J., Guo, Y.X., Zhou, B., Tong, X.: Indoor scene generation from a collection of semantic-segmented depth images. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 15203–15212 (2021)
77. Yang, Y., Sun, F.Y., Weihs, L., VanderBilt, E., Herrasti, A., Han, W., Wu, J., Haber, N., Krishna, R., Liu, L., et al.: Holodeck: Language guided generation of 3d embodied ai environments. In: The IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR 2024). vol. 30, pp. 20–25. IEEE/CVF (2024)
78. Yeshwanth, C., Liu, Y.C., Nießner, M., Dai, A.: Scannet++: A high-fidelity dataset of 3d indoor scenes. In: Proceedings of the International Conference on Computer Vision (ICCV) (2023)
79. Yu, L.F., Yeung, S.K., Tang, C.K., Terzopoulos, D., Chan, T.F., Osher, S.J.: Make it home: automatic optimization of furniture arrangement. ACM Transactions on Graphics (TOG)-Proceedings of ACM SIGGRAPH 2011, v. 30,(4), July 2011, article no. 86 **30**(4) (2011)
80. Zhai, G., Örnek, E.P., Wu, S.C., Di, Y., Tombari, F., Navab, N., Busam, B.: Commonscenes: Generating commonsense 3d indoor scenes with scene graphs. arXiv preprint arXiv:2305.16283 (2023)
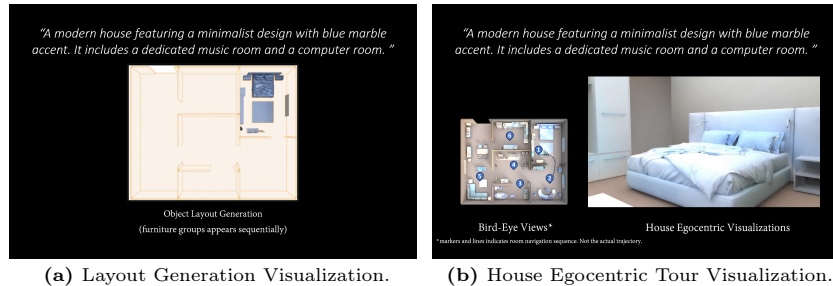
81. Zhang, Z., Yang, Z., Ma, C., Luo, L., Huth, A., Vouga, E., Huang, Q.: Deep generative modeling for scene synthesis via hybrid representations. ACM Transactions on Graphics (TOG) **39**(2), 1–21 (2020)
82. Zhao, W.X., Zhou, K., Li, J., Tang, T., Wang, X., Hou, Y., Min, Y., Zhang, B., Zhang, J., Dong, Z., et al.: A survey of large language models. arXiv preprint arXiv:2303.18223 (2023)

# Supplementary Material for AnyHome

This is the Supplementary Material for **AnyHome**. It includes additional visualizations and experimental results, featuring video visualizations that navigate through the interiors of generated houses (Sec. 6), trajectory generation visualizations (Sec. 7), layout refinement visualizations (Sec. 8), complex texture generation (Sec. 9), diverse floorplans (Sec. 10), and a qualitative and quantitative comparison with contemporary work (Sec. 12), highlighting our method's uniqueness. We present the time consumption breakdown for each component (Sec. 11). We also analyze the limitations of the current method (Sec. 13). Furthermore, we provide details on our method's settings, including implementation details (Sec. 14), utilized prompts (Sec. 15), and the rules and algorithms for object placement (Sec. 16).
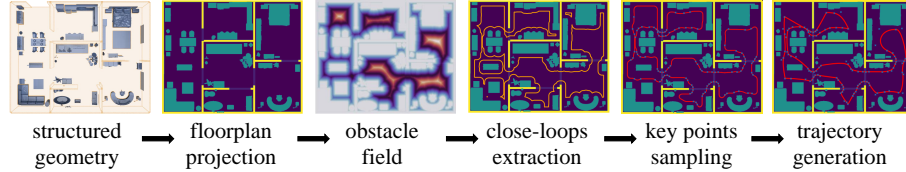
## 6    House Egocentric Visualizations with Videos

A Supplementary Video is included to offer additional qualitative results. Fig. 9 presents screenshots taken from the Supplementary Video. This video demonstrates the sequential placement of furniture groups within the layout generation process (see Fig. 9a), effectively showcasing the nuances of our placement algorithm (refer to Algorithm 1). Furthermore, it features egocentric navigation through the generated rooms (see Fig. 9b), highlighting how our methods for egocentric refinement and in-painting achieve coherent layouts and create realistic textures across different viewpoints.



(a) Layout Generation Visualization.          (b) House Egocentric Tour Visualization.

**Fig. 9:** Supplementary video provide layout generation visualization and house egocentric visualizations.

## 7    Egocentric Trajectory Generation



structured geometry → floorplan projection → obstacle field → close-loops extraction → key points sampling → trajectory generation
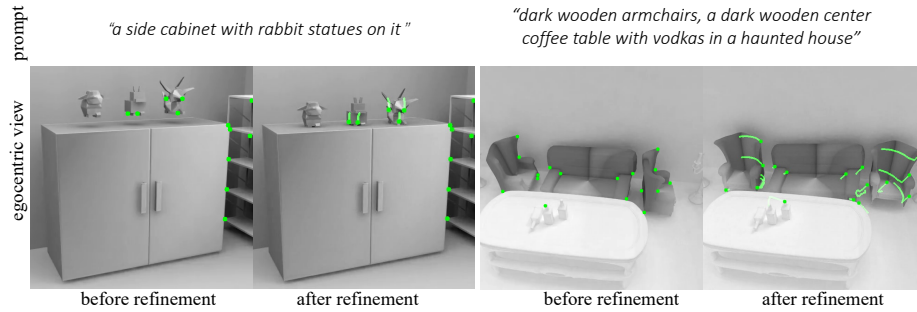
**Fig. 10: Egocentric Trajectory Generation Pipeline.** Given the base mesh generated from amodal representations, floorplan projection, obstacle field and close-loops are extracted sequentially. Key points are then sampled and interpolated to obtain the final trajectory.

Fig. 10 delineates the process of egocentric trajectory generation, an essential element that fosters object-level layout refinement and texture in-painting. Starting with the structured geometry derived from amodal representations, we project this geometry onto a 2D plane to obtain a floor plan. From this floor plan, we extract an obstacle field, which calculates the nearest distance from each point to furniture or walls. Utilizing these distance fields, we identify closed loops that are subsequently vectorized for refinement. This refinement process includes the elimination of short loops and areas of high curvature, along with the application of Gaussian smoothing to address minor irregularities. To facilitate visibility of objects and ensure navigation through doorways during exploration, keypoints are meticulously placed around these critical features on the closed loops. The final trajectory is constructed by interpolating between these strategically sampled keypoints, ensuring a coherent path that enhances the scene's navigability and visibility.
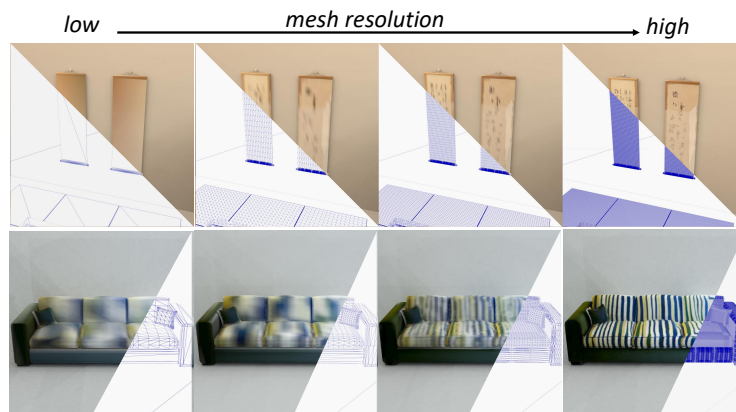
# 8    Layout Refinement Visualization

Fig. 11 visualizes object layouts before and after refinement. The green high-lighted lines illustrate the objects' positions throughout the optimization process from a single viewpoint. The key-points are extracted by Scale-Invariant Feature Transform (SIFT) and tracked by optical flow. The results indicate that the refinement process effectively anchors floating items (left) and improves object orientation while correcting asymmetry (right).



**Fig. 11: Layout Refinement Visualization.** The refinement trajectory is visualized in green color through optical flow. Scene layout shows enhanced alignment with the textual prompt following layout refinement.

## 9   Complex Texture

Fig. 12 shows Anyhome's ability to generate both complex details (the calligraphy) and decorative patterns (the stripe). The complexity of the texture is limited by the mesh resolution of object datasets (Objaverse and 3D Future)—complex patterns cannot be applied to a surface consisting of only two triangular meshes. Instead, we can re-mesh these objects in the original dataset to a higher resolution to accommodate more complex textures, albeit at the cost of increased inference time. As shown in the figure, the clarity of both the calligraphy and stripe patterns improves with higher mesh resolution.
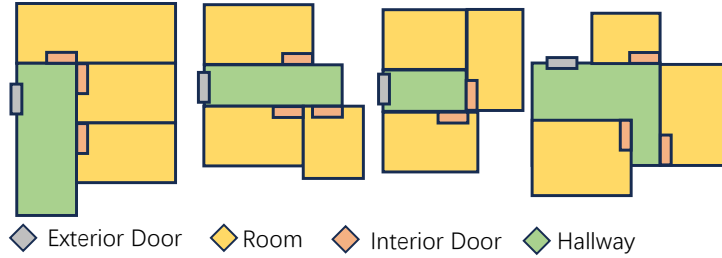


**Fig. 12: Generating Complex Texture.** AnyHome is able to generate complex texture as mesh resolution increases from left to right.
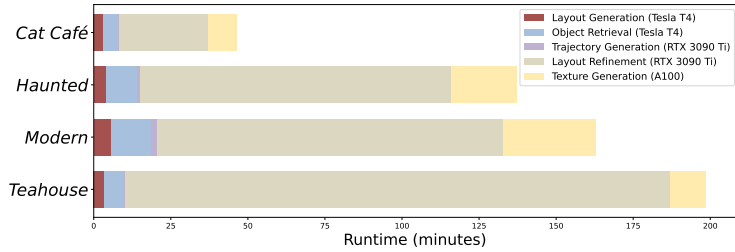
## 10    Floorplan Diversity

Fig. 13 shows that AnyHome can generate diverse floorplans given the same prompt `"A 1B1B apartment with a hallway."`. The generation of novel room types is made possible by the LLM's ability of mapping unconventional room type (hallway) to types that exist in the model's training set (i.e., entrance). At the same time, the diversity results from HouseGAN++'s generalizability.



◇ Exterior Door     ◇ Room     ◇ Interior Door     ◇ Hallway

**Fig. 13: Generating Novel Room-Type.** Given the prompt `"A 1B1B apartment with a hallway."`, AnyHome is able to generate novel room types (i.e., the hallway) and diverse floorplans.
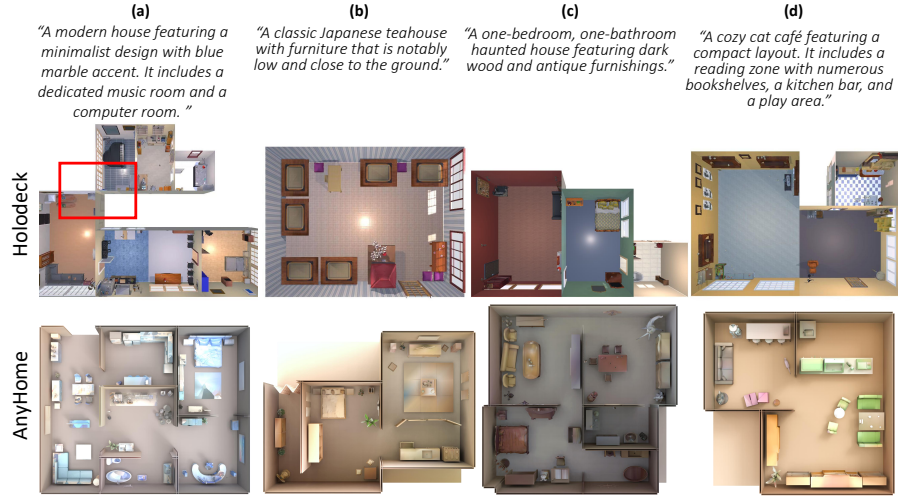
## 11    Time Consumption

Fig. 14 presents a detailed breakdown of time consumption for the four sample houses included in the main paper.



**Fig. 14: Time Consumption.** The time consumption of each component for examples shown in the main paper.

## 12    Comparison with Contemporary Works



**Fig. 15: Qualitative Comparison with Holodeck. AnyHome** proves its uniqueness and competitiveness through its ability to generate more logical and complex floorplans, achieve a broader array of layout designs, and provide consistent and customizable texture creation when compared to contemporary works such as Holodeck.

Holodeck [77] represents a contemporary method for text-conditioned indoor scene generation. Although direct comparisons with contemporaneous works are not obligatory, our goal is to highlight **AnyHome**'s unique attributes and competitive edge. Fig. 15 offers a qualitative comparison with Holodeck, illustrating differences in several key aspects.

**Floorplan Generation.** Holodeck is capable of generating regular floorplans directly through LLMs. However, it often produces impractical floorplans for complex scenes, as highlighted in red in panel **a**, and the floorplans tend to be overly simplistic, as observed in panels **b** and **c**. In contrast, **AnyHome** generates realistic and diverse floorplans, featuring various functional rooms that enhance the overall functionality of the house.

**Room Object Selection.** While Holodeck selects appropriate furniture for common room types, such as a bed for a bedroom in panel **c**, it struggles to identify suitable objects for less common room types, leading to empty spaces like the cat café room in panel **d**. It also shows limitations in the quantity and diversity of objects, evidenced by the sparse furnishings in panels **c** and **d** and repetitive furniture in panel **b**. In contrast, **AnyHome** adeptly selects furniture corresponding to both the room type and the house's description, resulting in efficiently utilized and densely furnished spaces.

**Table 2:** Quantitative comparison between Holodeck and **AnyHome** using GPT-4V.

| Method | Prompt-Align | Layout | Object | Texture | Overall |
|---|---|---|---|---|---|
| Holodeck [77] | 7.8 | 6.1 | 4.3 | 5.8 | 6.0 |
| **AnyHome** | **9.0** | **8.7** | **7.8** | **6.0** | **7.8** |

**Room Layout Generation.** Holodeck tends to position objects along walls, resulting in less varied layouts. Conversely, **AnyHome** facilitates more realistic and varied room layouts, thanks to placement rules informed by Large Language Models (LLMs).

**Object Texture.** Objects in Holodeck's scenes exhibit a wide range of textures, inherently tied to the object meshes, leaving limited room for customization. On the contrary, **AnyHome** applies texture to object meshes through diffusion models with style-specific prompts, enabling texture customization and visual domain adaptation. This approach produces scenes with enhanced texture consistency and better alignment with the input description, exemplified by the creation of a "blue marble accent" in panel **a**.

As Holodeck [77] does not support egocentric views or prompts, a direct comparison with **AnyHome** using the metrics outlined in the main paper is not feasible. Consequently, following the approach of [73], we employ GPT-4V for evaluation by providing it with top-down views of the generated scenes. Table 2 compares our method with Holodeck across various dimensions: the scene's alignment with the given prompt (*Prompt-Align*), the realism of house plans and room layouts (*Layout*), the authenticity of objects (*Object*), and the consistency of textures (*Texture*). The specific prompts used for these evaluations and the scoring scale are detailed in Fig. 16.

Quantitative analysis reveals that **AnyHome** excels in aligning scenes with the provided prompts, achieving an impressive average score of 9.0 out of 10 for *Prompt-Align*. This superiority stems from its capability to generate realistic floorplans for any designated room types, accurately retrieve objects corresponding to a room's function and style, and apply texture in line with the prompt's specifications. In individual assessments, **AnyHome** also stands out. Its two-stage approach for floorplan and room layout generation yields average scores of 8.7 for *Layout* and 7.8 for *Object*, demonstrating the generated geometry's plausibility and the robustness of this generation. Notably, Holodeck scores relatively low in *Object*, primarily due to its tendency to create sparsely furnished rooms, as depicted in Fig. 15. Both **AnyHome** and Holodeck show limitations in achieving texture consistency (*Texture*), with **AnyHome** occasionally painting a single object in inconsistent colors, as illustrated in Fig. 17. This issue arises from inconsistencies in viewpoint and the images generated from the multi-view diffusion inpainting model [65]. Overall, **AnyHome** attains a score of 7.8, indicating the presence of flaws but affirming the method's practical applicability in real-world scenarios.

Our task here is to rate a 3D indoor scene generated from a text description based on its top-down image. Your rating should be obtained based on the following instructions, and you should rate on a scale of 10.

# Instruction

**Prompt Alignment.** Focus on how well the scene corresponds to the given text description. An ideal scene should accurately reflect all the layout structure and in-painting style mentioned in the text prompt, capturing the corresponding attributes as described. Please first describe the layout and style of each scene and then evaluate how well it adheres to the requirements in the input description.

**Layout Plausibility.** Look at both the floor plan (each room's size and arrangement in the house) and furniture arrangement. Pay attention to any layout that is unrealistic and messy, for example, disconnected rooms or chairs far away from the table. An ideal scene should have a regular floor plan in adequate sizes and arrangements, with furniture laid suitably in groups and in the right position.

**Object Plausibility.** Focus on the selection of objects in each room. Deduce whether these objects are sensible to be placed in the house as per the text description. You should also evaluate whether the objects in this scene are diverse in terms of quantity and style.

**Texture Consistency.** Look over the texture of each furniture in each room, and check whether they are consistent with each other. Scenes that have objects showing different textures in multi-view perspectives should be penalized in this criteria. The objects' texture should also echo that of the room and the floor.

Rate the scene on each degree carefully on a scale of 10. A score of 10 means that the scene is perfect and the generated scene is in a way that matches the best design in the world. A score of 8 means that there are some flaws in the scene, but it is still sensible and suitable for real-world applications. A score of 6 indicates that the flaws are obvious in the real-world applications. A score of 4 indicates that the criteria are somewhat met but still plausible. A score of 2 indicates that the criteria are hardly met. A score of 0 indicates a complete mismatch to the criteria. After that, please obtain a total final rating for the scene, considering the scene's performance on all the degrees.

# Output format

Before rating the scene on each degree, please provide a short analysis of its performance on each of the above-mentioned evaluation criteria. The analysis should be concise and accurate. You should format your response as below:
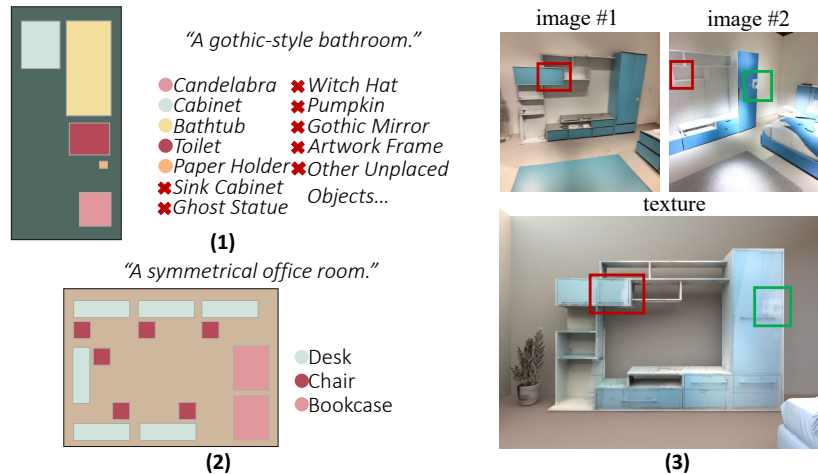
1. Text prompt and Asset Alignment: The scene... Degree Score: x/10.
2. Layout Plausibility: The scene... Degree Score: x/10.
3. Object Plausibility: The scene... Degree Score: x/10.
4. Texture Consistency: The scene... Degree Score: x/10.

Total Score: x/10

**Fig. 16:** Prompts for GPT-V4 evaluation.

## 13    Limitations and Failure Cases

**AnyHome**, as a pioneering approach to generating customized, house-scale scenes, faces several challenges. First, despite repeatedly emphasizing room dimensions to the Large Language Model (LLM), it occasionally places too many objects in limited spaces (see Fig. 17 **(1)**), highlighting its limitations in grasping sizes and spatial relationships. Although the placement algorithm stops adding objects once space is exhausted, ensuring the final layout remains viable, the initial list of objects and the suggested placement rules from the LLM do not always result in a logically arranged layout. Second, while **AnyHome** includes rules capable of generating symmetrical layouts—such as classrooms with tables and chairs in orderly rows and columns, as shown in Fig. 7 of the main paper—the LLM sometimes misunderstands these rules. It selects rules that place furniture reasonably but fails to comply with the style specified in the prompt, as evident in (Fig.17 **(2)**), where desks and chairs are realistically but incorrectly arranged around the walls. Lastly, our egocentric inpainting approach, though it facilitates realistic texture customization, encounters issues with view consistency. Fig. 17 **(3)** illustrates these challenges, where red boxes highlight asymmetrical textures resulting from the disparate images generated by the multi-view diffusion inpainting model [65] across different viewpoints, and green boxes show the loss of texture detail due to viewpoint inconsistencies.



**Fig. 17: Limitations.** (1) Generation of overly extensive lists by the LLM for confined spaces. (2) Inability to adhere to particular layout constraints stemming from the LLM's interpretation of the rules. (3) Asymmetrical textures arising from inconsistencies in input images sourced from various viewpoints.

## 14  Implementation Details

**LLM-guided Layout Generation.** We utilize OpenAI's `gpt-4-1106-preview` model as the LLM for our system and set the `temperature` to `0.7` to allow diversity in each generation. Simultaneously, we set `Top-p` to `1` and apply no penalty to the next token prediction.

**Diffusion-guided Refinement.** For layout refinement, we employ Stable Diffusion 2 as our diffusion model. In the diffusion process, we select a batch of 8 images representing multiple views, all associated with the same prompt, for optimization. Beyond using normal maps paired with diffusion prompts, we also introduce a technique where objects are rendered in random colors—such as "red," "blue," "green," "yellow," "purple," and "orange"—and these images are paired with prompts modified to reflect these colors for enhanced refinement. For instance, if rabbits are rendered in blue and a cabinet in red, the original ego-centric diffusion prompt "a side cabinet with rabbit statues on it" is adjusted to "a red side cabinet with blue rabbit statues on it." This color-modified approach allows for more precise refinement by providing the diffusion model with additional context about the scene's color composition, thus improving the alignment between the generated images and the specified layout and appearance.

**Diffusion-guided Inpainting.** We use MVDiffusion as our multi-view, depth-conditioned inpainting model [65], setting a `guidance_scale` of `15` to increase the influence of text prompts in generation and an `overlap_filter` of `0.1` to ensure maximal consistency.

## 15  Prompts

Fig. 18, 19, 20, 21, 22 and 23 present the complete prompt templates employed by **AnyHome**. The prompt $p_{\text{floorplan}}$ is designed to generate a bubble diagram for house floorplans, while $p_{\text{map}}$ translates unconventional room types into those included in the RPLAN dataset [74] to facilitate floorplan generation. These prompts are combined for efficient generation. The $p_{\text{room}}$ prompt produces a room constraint graph, detailing the description, dimensions, and placement of objects. It is divided into two distinct prompts, $p_{\text{furniture}}$ and $p_{\text{ornament}}$ for furniture and ornaments respectively, to address the Large Language Model (LLM)'s limitations in generating diverse scenes with complex descriptions in a single prompt round. $p_{\text{furniture\_edit}}$ and $p_{\text{ornamen\_edit}}$ are the corresponding prompts for editing. The $p_{\text{appearance}}$ prompt directs the LLM to provide inpainting prompt for each depth map by describing the objects and their image bounding boxes, which are calculated from viewpoints sampled by the egocentric trajectory.

Our prompt templates are designed to: (1) engage the LLM as an assistant for 3D scene generation, (2) outline requirements systematically, and (3) specify the desired output format, specifically requesting responses in JSON to streamline parsing and reduce irrelevant information. To improve adherence to expected outputs, we include task examples for the LLM. Notably, to address instances where LLMs overlook specified room types or styles—such as placing a double

bed in a living room or choosing modern furniture for a medieval castle—we reinforce these details in the prompts to enhance consistency with given inputs.

**Task:** You are a talented Architectural Planner tasked with envisioning the floorplan for a house described as <house_description>.

**Requirements:**
1. **Complete Room Inventory:** Compile a comprehensive list of all rooms within the house. If a room type appears multiple times, append an index to differentiate them. Aim for a diverse assortment of rooms, covering a wide range of functionalities. For example, if the house features two bedrooms and one dining room, list them as: "[bedroom1, bedroom2, dining room1]".
2. **Standardized Room List:** Adapt the previously listed rooms into a standardized set of room types based on their functionality and the house's style. Use only the following room types: *kitchen, storage, bathroom, study_room, balcony, living_room, bedroom, entrance, dining_room, and unknown*. Transform unique rooms to the closest match from the predefined list, appending an index if necessary. If a room does not closely match any predefined type, label it as "unknown".
3. **Room Connections:** Map out the connectivity between rooms, detailing which rooms are directly accessible from each other. Present this information as a list of tuples, indicating room pairs that share a connection. For instance, if the dining room connects to both bedroom1 and bedroom2, but there is no direct connection between the two bedrooms, list the connections as: "[[dining_room1, bedroom1], [dining_room1, bedroom2]]". The room names should be from the complete room list that you've generated at requirement 1.
4. **Front Door Locations:** Identify the rooms that house the main entrances to the dwelling. Specify each room that contains a front door, considering it as the primary access point to the house.

**Output:** Provide the information in a valid JSON structure with no spaces.
{
   "complete_room_list": [...],
   "modified_room_list": [...],
   "connections": [...],
   "front_door": [...]
}

**Fig. 18:** $p_{\text{floorplan}}$ and $p_{\text{map}}$ for floor plan generation.

**Task**: You are an awesome 3D Scene Designer. Design a 3D indoor scene for a <room_type> located within a <house_description>. Ensure that the design fits within an area of <room_area> square meters. Provide the details as a scene graph, structured in the JSON format.

**Requirements:**
**1. Furniture List:**
 - Enumerate the furniture in the <room_type>. If an item appears multiple times, append an index. Remember that the furniture should be suitable for a <house_description> and maintain the style of the house. For example, you shouldn't only generate two sets of tables and chairs for an office room.
 - Use only from this predefined list: *children_cabinet, nightstand, bookcase, wardrobe, coffee_table, corner_table, side_cabinet, wine_cabinet, tv_stand, drawer_chest, shelf, round_end_table, double_bed, queen_bed, king_bed, bunk_bed, bed_frame, single_bed, kids_bed, dining_chair, lounge_chair, office_chair, dressing_chair, classic_chinese_chair, barstool, dressing_table, dining_table, desk, three_seat_sofa, armchair, loveseat_sofa, l_shaped_sofa, lazy_sofa, chaise_longue_sofa, stool, kitchen_cabinet, toilet, bathtub, sink.*
 - You need to generate as many furnitures as you can.
 - Example: If there are two lazy sofas and an armchair in the room, list them as: "[lazy_sofa1, lazy_sofa2, armchair1]".
**2. Furniture Descriptions:**
 - For each furniture, provide a description of its aesthetic shape and structure considering the house and room's styles.
 - Example: "A wood-carved three-seated sofa with two armrests at the sides."
 - Return format: {<furniture_name>: <description>}
**3. Furniture Sizes:**
 - For each furniture, provide its dimensions (length, width, height) in meters, even if the dimensions repeat with other furnitures, keeping in mind the <room_area> square meters room area.
 - Return format: {<furniture_name>: [length(meters), width(meters), height(meters)]}
**4. Furniture Groups & Placement Rules:**
 - Separate the furnitures into groups that will exist together in the scene, such as a bed and two nightstands. For those furnitures that are unrelated with each other, put them in different groups. Put the important furniture groups at first when returning your answer. For example, put the funiture group containing the sofa, tv stand and coffee table at first in a living room.
 - For each group, you should determine only one anchor furniture. For example, the anchor for a bed and two nightstands is the bed.
 - Then, you should decide on the placement rule for this anchor furniture. You can use only the following anchor rules:
  (1) "place_center" which places the anchor furniture at the center of the room.
  (2) "place_next_wall" which places the anchor with its side against a segment of the wall, useful for rooms with rows and columns of furniture like a classroom.
  (3) "place_wall" which places the anchor with its back against a segment of the wall.
  (4) "place_corner" which places the anchor at a corner.
  (5) "place_next(another_anchor_name, x)" which places the anchor furniture beside another anchor specified in the parameter with a buffer distance of x meters. The "another_anchor_name" is the another anchor's name in the furniture list.
 - Last, for the other non-anchor furnitures in the group, you can only use the following non-anchor rules:
  (1) "place_front(x)" which places the furniture in front of the anchor with a buffer distance of x meters, like a TV stand before a coffee table.
  (2) "place_beside(x)" which places the furniture beside the anchor with a buffer distance of x meters, like a nightstand beside a bed.
  (3) "place_around(x)" which places the furniture around the anchor with a buffer distance of x meters, like placing four chairs around a dining table.
 - You can set x to 0 if you want the furniture to be placed right next to the anchor.
 - CAUTIOUS: Anchor furnitures can only use the anchor rules, while non-anchor furnitures can only use the non-anchor rules. For example, DO NOT use "place_next(another_anchor_name, x)" for non-anchor furnitures, use "place_beside(x) instead".
 - Return format: [[[<anchor_furniture_name>, <anchor_placement_rule>], [<non-anchor_furniture_name>, <non-anchor_placement_rule>], [<non-anchor_furniture_name>, <non-anchor_placement_rule>], ...], ...]

**Output**: Provide the information in a valid JSON structure with no spaces. I'll give you 100 bucks if you help me design a perfect scene and return it in the right format:

```
{
  "complete_room_list": [...],
  "furniture_descriptions": {...},
  "furniture_sizes": {...},
  "furniture_groups_and_placement_rules": [...]
}
```

**Fig. 19:** $p_{\text{room}}$ for room layout generation - the $p_{\text{furniture}}$ part.

**Task:** You are an awesome 3D Scene Designer. Design the ornaments added to an existing <room_type> located within a <house_description> that has an area of <room_area> square meters. These are the furnitures in the room right now: <existing_furniture_list>.

**Requirements:**

**1. Ornament List:**
  - Enumerate all the ornaments in this room. Any objects beside *children_cabinet, nightstand, bookcase, wardrobe, coffee_table, corner_table, side_cabinet, wine_cabinet, tv_stand, drawer_chest, shelf, round_end_table, double_bed, queen_bed, king_bed, bunk_bed, bed_frame, single_bed, kids_bed, dining_chair, lounge_chair, office_chair, dressing_chair, classic_chinese_chair, barstool, dressing_table, dining_table, desk, three_seat_sofa, armchair, loveseat_sofa, l_shaped_sofa, lazy_sofa, chaise_longue_sofa, stool, kitchen_cabinet, toilet, bathtub, and sink* are considered ornaments.
  - You can generate any ornament. Examples might include a knight statue, a witch's hat, and a Gothic clock.
  - Please take into consideration the area and style of the room when you are generating.
  - Generate as many ornaments as you can, and be CREATIVE with the ornaments.
  - Return format: [<ornament1>, <ornament2>, ...]

**2. Ornament Descriptions:**
  - For each ornament, provide a description of its aesthetic shape and structure considering the house and room's styles.
  - Example: "A classic witch hat, made of black velvet with a wide brim."
  - Return format: {<ornament_name>: <description>}

**3. Ornament Sizes:**
  - For each ornament, provide its dimensions (length, width, and height) in meters, keeping in mind the <room_area> square meters room area.
  - Return format: {<furniture_name>: [length(meters), width(meters), height(meters)]}

**4. Ornament Placements:**
  - For each ornament, provide its placement rule in the scene; below are the rules that you can use:
  (1) "place_center" which places the ornament at the center of available spaces in the room.
  (2) "place_next_wall" which places the ornament with its side against a segment of the wall.
  (3) "place_wall" which places the ornament with its back against a segment of the wall.
  (4) "place_corner" which places the ornament at a corner.
  (5) "place_front(x, anchor)" which places the ornament in front of an existing anchor furniture in the given list above with a buffer distance of x meters, like placing a mirror in front of a cabinet.
  (6) "place_beside(x, anchor)" which places the ornament beside an existing anchor furniture in the given list above with a buffer distance of x meters, like placing a vase beside a TV stand.
  (7) "place_around(x, anchor)" which places the ornament around an existing anchor furniture in the given list above with a buffer distance of x meters, like placing a set of candelabras around a dining table.
  (8) "place_top(x, anchor)" which places the ornament on top of an existing anchor furniture in the given list above with a buffer distance of x meters, like placing a vodka bottle on a table.
  (9) "place_on_wall(x)" which places the ornament on the wall at a height of x meters, like placing a clock on the wall.
  - For example, if a pear is placed 1 meter beside "table1", your response should be {"pear", "place_beside(1, table1)"}.
  - You can set x to 0 if you want the buffer distance to be zero.
  - Return format: {<ornament_name>: <placement_rule>}

**Output:** Provide the information in a valid JSON structure with no spaces. I'll give you 100 bucks if you help me design a perfect scene and return it in the right format:

```
{
  "ornament_list": [...],
  "ornament_descriptions": {...},
  "ornament_sizes": {...},
  "ornament_placements": [...]
}
```

**Fig. 20:** $p_{\text{room}}$ for room layout generation - the $p_{\text{ornament}}$ part.

**Task:** You are an awesome 3D Scene Designer. Your task is to create a prompt for an image generation Diffusion model based on a description of an indoor scene captured in a 256x256 image. Your prompt should reflect the scene's overall style, include all objects present in the image, and describe their aesthetics and spatial relationships based on the provided details.

**Image Description Details:**
**1. Scene's Overall Style:**
  - Below is the overall style of the scene. Incorporate the given style into your final prompt to ensure consistency with the scene's ambiance.
  - Style: "<house_description>"
**2. Object List:**
  - Below is the list enumerating all the objects in the image. Include all objects in the image in your final prompt.
  - Object List: <object_list>
**3. Object Descriptions:**
  - Below are the descriptions of each object's aesthetic shape and structure. Simplify each object's description to include its color (texture) and a brief description along with the object's name.
  - Use concise descriptions like "a dark wooden double bed" or "a stone sword with Japanese carvings."
  - Object Descriptions: "<object_descriptions>"
**4. Object Image Bounding Boxes:**
  - Below are the 2D image bounding boxes for the objects in this 256x256 image, with the x coordinate starting from the top-left corner going right, and the y coordinate starting from the top-left corner going down. Use the 2D image bounding boxes to describe the spatial relationships between objects. Bounding boxes are given as [(x1, y1), (x2, y2)], where (x1, y1) is the top-left corner, and (x2, y2) is the bottom-right corner. Highlight relationships such as "on top of" and "at the back of," and list objects without specifying side-by-side relationships.
  - Object Bounding Boxes: "<object_bounding_boxes>"

**Output Requirements:** Craft the final prompt for the image generation model, combining the objects' descriptions with their spatial relationships. Focus on listing all objects and their relationships, separated by commas, without concerning for grammatical correctness. An example might be "a green pear on top of a white modernistic table with a smooth texture, a grey fridge which is closed," depicting a scene with a pear, a table, and a fridge.

**Fig. 21:** $p_{\mathrm{appearance}}$ for in-painting prompt generation.

**Task**: You are an awesome 3D Scene Designer. You are given a 3D indoor scene graph of a <room_type> located within a <house_description>. The scene graph is generated with the four requirements below:

**Requirements:**
**1. Furniture List:** ... (The same as $p_{furniture}$ above)
**2. Furniture Descriptions:** ... (The same as $p_{furniture}$ above)
**3. Furniture Sizes:** ... (The same as $p_{furniture}$ above)
**4. Furniture Groups & Placement Rules:** ... (The same as $p_{furniture}$ above)

Generated Data of the Scene Graph of the <room_type>:
<The latest scene graph transformed as a string>

Now, the customer wants to edit the scene graph. They have provided the following description of the changes they want to make:
<The editing prompt input by the user>

**Task:** Could you please update the floorplan based on the customer's request and provide the updated data strictly following the same requirements above? If the customer's request contrasts with the requirement above, please follow the requirements above. For example, if the customer wants to use "place_around(x)" as the placement rule for an anchor furniture, you should not follow it and should generate the placement rule based on the requirements above.

**Output**: Provide the information in a valid JSON structure with no spaces. I'll give you 100 bucks if you help me design a perfect scene and return it in the right format:

```
{
  "complete_room_list": [...],
  "furniture_descriptions": {...},
  "furniture_sizes": {...},
  "furniture_groups_and_placement_rules": [...]
}
```

**Fig. 22:** $p_{\mathrm{furniture\_edit}}$ editing furniture-related properties.

**Task**: You are an awesome 3D Scene Designer. You are given a 3D indoor ornament graph of a <room_type> located within a <house_description>. The ornament graph is a complement to the existing furnitures in the room, and these furnitures are <existing_furniture_list>. The ornament graph is generated with the four requirements below:

**Requirements:**
**1. Ornament List:** ... (The same as $p_{ornament}$ above)
**2. Ornament Descriptions:** ... (The same as $p_{ornament}$ above)
**3. Ornament Sizes:** ... (The same as $p_{ornament}$ above)
**4. Ornament Placements:** ... (The same as $p_{ornament}$ above)

Generated Data of the Ornament Graph of the <room_type>:
<The latest ornament graph transformed as a string>

Now, the customer wants to edit the ornament graph. They have provided the following description of the changes they want to make:
<The editing prompt input by the user>

**Task:** Could you please update the ornament graph based on the customer's request and provide the updated data strictly following the same requirements above?

**Output**: Provide the information in a valid JSON structure with no spaces. I'll give you 100 bucks if you help me design a perfect ornament graph and return it in the right format:

```
{
    "complete_room_list": [...],
    "furniture_descriptions": {...},
    "furniture_sizes": {...},
    "furniture_groups_and_placement_rules": [...]
}
```

**Fig. 23:** $p_{\text{ornament\_edit}}$ for ornament-related editing.

# 16 Placement Rules and Algorithms

Table 3 outlines the placement rules employed by **AnyHome** to govern the positioning of objects within a room. In addition to the parameter `x`, which specifies the buffer distance in meters, these rules also consider the dimensions of the object, the dimensions and positions of the anchor object (for rules like `place_beside` that reference an anchor), and a mask of the room indicating its boundaries and spaces not occupied by previously placed objects. The LLM determines the appropriate rule for each object, along with these parameters.

Algorithm. 1 details the furniture placement algorithm. The placement algorithm iteratively positions furniture groups $C_i$, starting with the anchor object for each group. This anchor is decided by the LLM. The anchor's placement utilizes its rule to identify viable spaces within the room (e.g., every available corner for `place_corner`), placing the anchor where constraints are met. Should suitable spaces be unavailable, the furniture group is omitted. This approach is rational, as the LLM is instructed to prioritize the generation of the most significant furniture groups first (refer to Figure 19 for details), ensuring ample free space for their placement. Subsequently, the placement rule for each object within the group assesses all viable spaces that adhere to the constraints, such as positioning in front of the anchor, and discards any object for which no appropriate space is found. These rules not only guarantee that objects are placed in accordance with their relative positions within the room but also prevent objects from being positioned outside room boundaries or overlapping with one another.

---

**Algorithm 1** Furniture Placement Algorithm

---

**Input:** Room mask $M_i$, Subgraphs $C_i$, placement functions $P$
1:   $B_i \leftarrow empty$
2: **for** subgraph $c_i$ in $C_i$ **do**
3:     /* Identify and place the anchor furniture */
4:     $n_{c_i}^{anchor} \leftarrow$ identify_anchor($c_i$)
5:     $p \leftarrow P(n_{c_i}^{anchor})$
6:     $b_{c_i}^{anchor} \leftarrow p(M_i, \text{anchor})$
7:     $B_i \leftarrow B_i \cup b_{c_i}^{anchor}$
8:     $M_i \leftarrow$ update_mask($M_i, B_i$)
9:     /* Place the subsequent furniture items */
10:     **for** each furniture item $n_{c_i}^j$ in $c_i$ **do**
11:       $p \leftarrow P(n_{c_i}^j)$
12:       $B_i \leftarrow B_i \cup p(M_i, n_{c_i}^j, b_{c_i}^{anchor})$
13:       $M_i \leftarrow$ update_mask($M_i, B_i$)
14:     **end for**
15: **end for**
**Output:** Set of furniture bounding boxes $B_i$

---

| Rule | Description |
|---|---|
| place_corner | Place the object in an available corner of the room. |
| place_next_wall | Place the object next to an available wall section. |
| place_wall | Place the object against an available wall section, oriented towards the room center. |
| place_center | Place the object at the room center, or the center of the spare area that no objects have been placed above. |
| place_next | Place an anchor object next to an existing anchor object at the same orientation as that anchor. This is used for aligning groups of tables and chairs in rows and columns. |
| place_beside(x) | Place the object next to the anchor object with a buffer distance of x, such as a bed, at the same orientation as the anchor. This is used for object groups like a bed and two nightstands. |
| place_around(x) | Place the object at an available edge around the anchor object with a buffer distance of x, oriented towards the anchor. This is used for object groups like a dining table and chairs. |
| place_front(x) | Position the object in front of an anchor object with a buffer distance of x, oriented towards the anchor. This is used for object groups like a sofa and a TV stand. |
| place_top(x) | Place the object on top of an anchor object, with a buffer distance of x, at the center of the spare area on top of the anchor that no objects have been placed above. |
| place_on_wall(x) | Place the object on an available wall section with a height of x, orientated towards the room center. |

**Table 3:** The placement rules for room layout generation.