

# SCENEFOUNDRY: Generating Interactive Infinite 3D Worlds

ChunTeng Chen<sup>1</sup> YiChen Hsu<sup>2</sup> Yiwen Liu<sup>1</sup> WeiFang Sun<sup>3</sup>  
TsaiChing Ni<sup>1</sup> ChunYi Lee<sup>4</sup> Min Sun<sup>2</sup> YuanFu Yang<sup>1</sup>

<sup>1</sup>National Yang Ming Chiao Tung University <sup>2</sup>National Tsing Hua University

<sup>3</sup>NVIDIA AI Technology Center <sup>4</sup>National Taiwan University

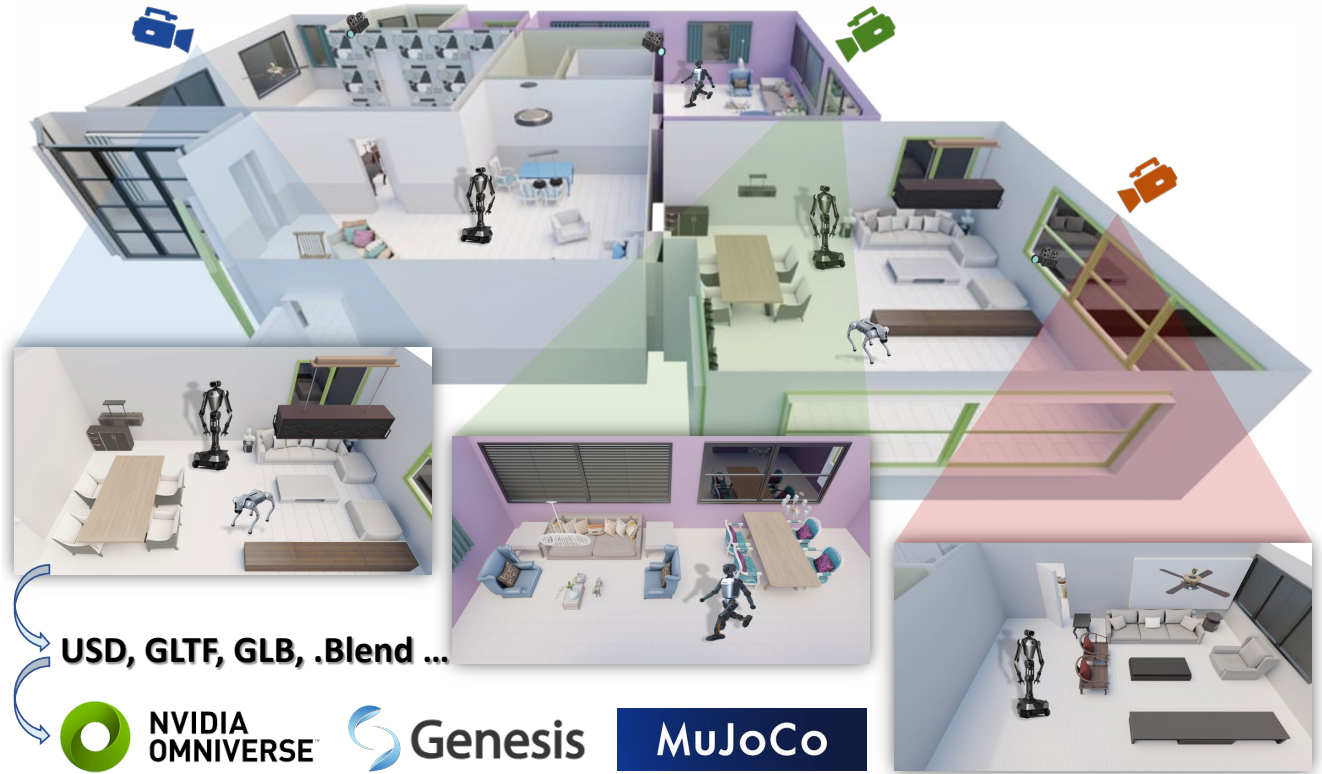


Figure 1. Overview of **SceneFoundry**. The framework generates apartment-scale 3D scenes from natural language prompts via LLM-guided floor plan generation, diffusion-based placement, and post-optimization ensuring articulated functionality and robot navigability.

## Abstract

The ability to automatically generate large-scale, interactive, and physically realistic 3D environments is crucial for advancing robotic learning and embodied intelligence. However, existing generative approaches often fail to capture the functional complexity of real-world interiors, particularly those containing articulated objects with movable parts essential for manipulation and navigation. This paper presents *SceneFoundry*, a language-guided diffusion framework that generates apartment-scale 3D worlds with functionally articulated furniture and semantically diverse lay-

outs for robotic training. From natural language prompts, an LLM module controls floor layout generation, while diffusion-based posterior sampling efficiently populates the scene with articulated assets from large-scale 3D repositories. To ensure physical usability, *SceneFoundry* employs differentiable guidance functions to regulate object quantity, prevent articulation collisions, and maintain sufficient walkable space for robotic navigation. Extensive experiments demonstrate that our framework generates structurally valid, semantically coherent, and functionally interactive environments across diverse scene types and conditions, enabling scalable embodied AI research.

## 1. Introduction

The ability to generate diverse, large-scale, and realistic 3D indoor environments is fundamental to the advancement of robotics [18], virtual reality, and embodied AI [14]. However, bridging the simulation-to-reality gap remains a significant hurdle [4], often because generated simulation environments lack the complexity, controllability, and functional realism of their physical counterparts.

Recent efforts [13, 15, 20] have focused on increasing the physical realism and interactability of simulated environments. While enhancing functional realism in specific aspects, such approaches can inadvertently compromise the visual realism or diversity of the layouts. A major limitation of many learning-based methods is their inability to generate complete apartment-scale layouts, as they often focus only on single rooms. Procedural generation frameworks like Infinigen [17] can produce such large-scale environments, but they are computationally intensive. Existing works [16, 19, 20] also often lack fine-grained control over crucial scene properties, which is essential for generating targeted training data distributions.

This paper introduces a multi-stage, controllable generative framework designed to generate apartment-scale 3D indoor scenes that are not only visually diverse but also semantically coherent and functionally sound. Our approach bridges the gap between high-level user intent, specified via natural language, and the generation of structurally valid layouts suitable for robotic interaction. We architect a pipeline that integrates semantic guidance with a suite of novel constraints to ensure that every generated scene is tailored to specific requirements and is physically usable. As shown in Figure 1, SceneFoundry generates photorealistic, apartment-scale, and controllable 3D scenes. Our contributions are summarized below.

- **LLM-based Parameter Space Guidance.** We introduce a module that translates abstract user commands into low-level parameters, enabling semantic control over the generative priors of a floor plan generator.
- **Novel Functional Guidance Mechanisms.** We introduce a set of differentiable constraints to enforce functional plausibility. This includes:
  - An **Object Quantity Control** for precisely enforcing the number of objects in a scene.
  - An **Articulated Object Collision Constraint** that penalizes configurations where functional parts are obstructed, ensuring interactability.
- **Walkable Area Control.** A final Walkable Area Control optimization is applied to the generated layout to refine spatial density and guarantee agent navigability.
- **Novel Evaluation Metrics.** To validate the effectiveness of the generation and control methods, we introduce new evaluation metrics that measure controllability.

## 2. Related Work

**Indoor Scene Layout Generation.** The automated generation of 3D indoor scenes is a long-standing challenge in computer graphics and vision. Early approaches relied on procedural generation, employing rule-based grammars or optimization techniques to synthesize layouts. A prominent recent example, Infinigen [17], utilizes procedural methods combined with simulated annealing to generate high-fidelity, apartment-scale layouts. While capable of producing complex and realistic results, these methods are often computationally intensive, time-consuming, and difficult to control without expert knowledge of the underlying rules.

Learning-based approaches are now dominant. Autoregressive models, such as ATISS [16], generate objects sequentially. While this models inter-object relationships, it suffers from error accumulation, slow sampling, and difficult holistic editing. Diffusion models [12] are a powerful alternative. Methods like DiffuScene [19] generate all object parameters in parallel, offering superior holistic coherence, editing flexibility, and state-of-the-art quality [6]. We therefore adopt this paradigm.

Effective generative modeling critically depends on both dataset quality and structural consistency. While large-scale datasets of real-world scans like ScanNet [5] and Matterport3D [3] are invaluable for reconstruction and navigation, they are less suitable for *generative* tasks. We therefore utilize clean, CAD-based datasets, 3D-FRONT [8] and GPartNet [9], which offer well-structured geometry and part-level semantics ideal for controllable 3D scene synthesis.

**Guidance of Diffusion Model.** Controlling generative models is crucial. Early studies introduced classifier guidance [11], which leverages the gradient of an external classifier to steer the sampling process. This gradient-based steering concept was later generalized. The core idea, often referred to as diffusion posterior sampling, is highly flexible and allows for guidance through any differentiable function, not just a classifier. This principle is ideal for enforcing functional 3D constraints, and recent work has begun to explore its use for physical plausibility or robot reachability [1, 11, 20]. We adopt this posterior sampling approach for its modularity, training a single unconditional model and applying diverse constraints at test time while avoiding the cost of multiple specialized conditional models.

## 3. SceneFoundry

Our framework adopts a multi-stage pipeline to generate controllable, apartment-scale 3D scenes for robot training, as illustrated in Figure 2. The pipeline begins with an LLM-based parameter space generation (Sec. 3.1) that translates user prompts into low-level parameters for floor plan gen-

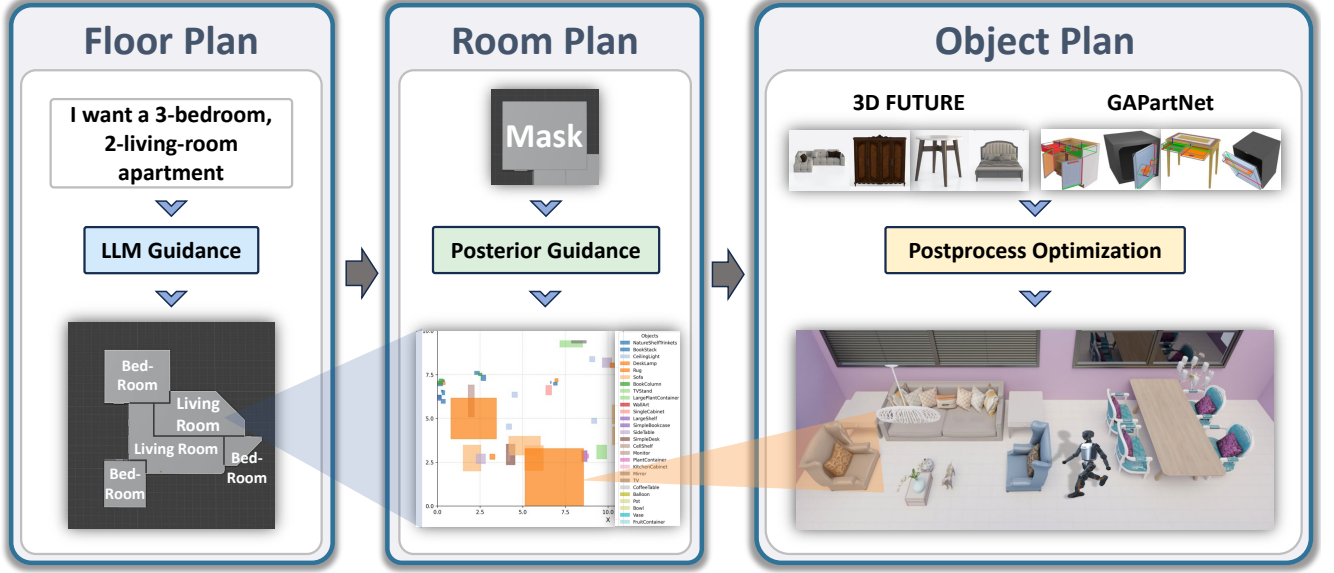


Figure 2. Overview of our apartment-scale generation pipeline. An LLM first guides procedural floor plan generation (Sec. 3.1), diffusion posterior guidance generates plausible room bounding boxes (Sec. 3.2, Sec. 3.3, Sec. 3.4), and 3D assets from 3D-FRONT/GAPartNet are refined via post-optimization to complete the layout (Sec. 3.5).

eration. A diffusion model employing posterior sampling (Sec. 3.2) then populates these layouts with furniture assets.

To ensure functional viability, we integrate three control mechanisms. During sampling, the model is guided by differentiable guidance functions: Object Quantity Control (Sec. 3.3) and the proposed Articulated Object Collision Constraint (Sec. 3.4) to maintain usability of movable parts. A Walkable Area Control post-processing step (Sec. 3.5) further refines spatial density to guarantee navigability. We

also introduce a set of evaluation metrics (Sec 3.6) to quantitatively evaluate the effectiveness of our methods.

### 3.1. LLM-Guided Parameter Space Generation

We adopt the Infinigen [17] framework, which generates room layouts through a simulated annealing process governed by twelve reward functions. However, its high-dimensional parameter space is not intuitive to control for users. To enhance usability, we design an LLM-based parameterization framework that interprets natural-language prompts and produces corresponding parameters for these functions, as shown in Figure 3. This semantic-to-parameter mapping converts abstract user descriptions into concrete floor plans while preserving the stochastic diversity of Infinigen.

### 3.2. Diffusion Posterior Sampling

Our method steers the reverse diffusion trajectory using a composite guidance function,  $\varphi(\cdot)$ , which enforces explicit constraints on the generated 3D scenes. This approach adapts the principles of diffusion posterior sampling to ensure structural validity in the generated layouts, as shown in the reverse process part of Figure 4.

**Object Feature.** Following [19], we represent a 3D scene  $x$  as an unordered set of  $N$  objects,  $\{\mathbf{o}_i\}_{i=1}^N$ . Each object  $\mathbf{o}_i$  is a vector  $[l_i, s_i, \theta_i, c_i, f_i]$  encoding its location, size, orientation, semantics, and a latent shape feature. This latent space is derived from a pre-trained VAE, following [19, 20]. The generated latent feature  $f_i$  is used for a nearest-neighbor search to retrieve the best-matching asset from ei-

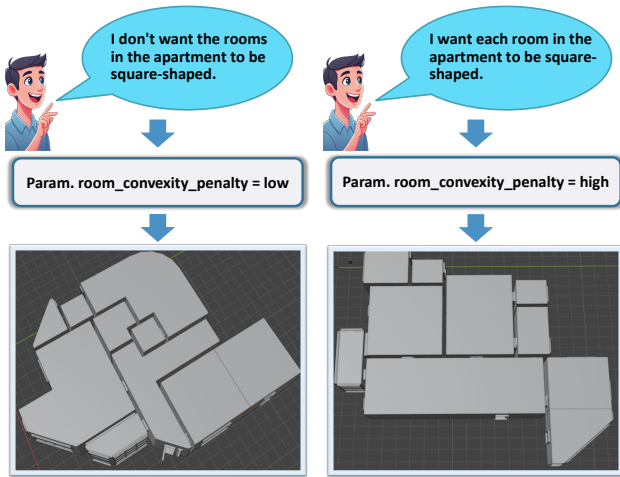


Figure 3. Illustration of our LLM-based Guidance. A low penalty (left) produces diverse, non-rectilinear layouts, whereas a high penalty (right) enforces square-shaped room layouts.

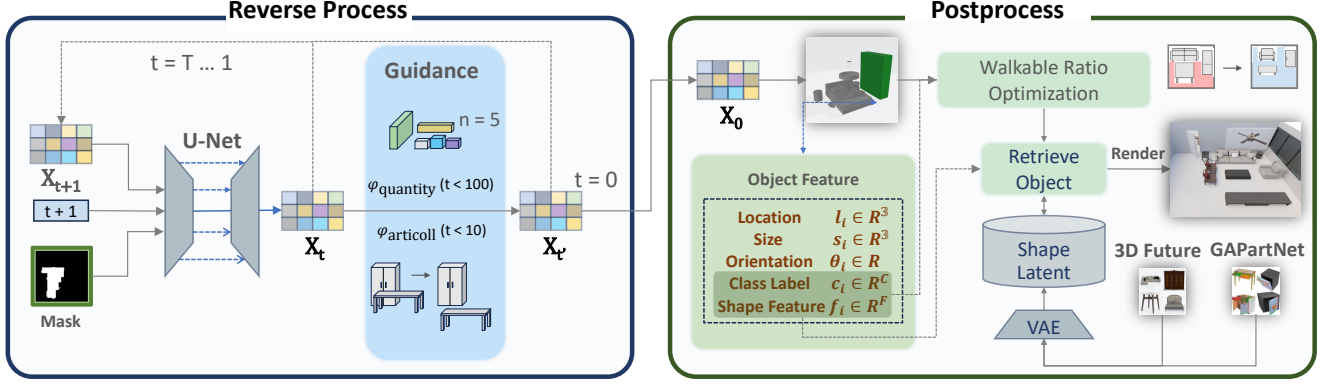


Figure 4. Guidance scheduling during the reverse diffusion process. Object quantity control is applied at  $t < 100$  and articulated collision constraint at  $t < 10$ , followed by a final walkable-ratio optimization at  $t = 0$  to generate a realistic scene.

ther 3D-FRONT or GAPartNet, as shown in postprocess part of Figure 4. This allows us to compose novel scenes that cohesively integrate both static and articulated objects.

**Training.** We adopt a *constraint-guided learning* strategy. Instead of a standard denoising objective, the model  $\epsilon_\theta$  is trained to predict the noise  $\epsilon$  while simultaneously anticipating the constraint gradient  $\mathbf{g}$  (from our constraint functions  $\varphi$ ). This is optimized by minimizing a guided  $\mathcal{L}_2$  loss:

$$\mathcal{L} = \mathbb{E}_{t, \mathbf{x}_0, \epsilon} [\|\epsilon - \lambda \Sigma \mathbf{g} - \epsilon_\theta(\mathbf{x}_t, t, \mathcal{F})\|_2^2] \quad (1)$$

This approach embeds knowledge of the constraints directly into the model weights during the training phase.

**Sampling.** During inference, we perform an iterative reverse process starting from  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . At each step  $t$ , the model first predicts the parameters  $(\mu_\theta, \Sigma_\theta)$  of the unguided posterior  $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$ . We then compute the gradient of our composite guidance function,  $\nabla_{\mathbf{x}_t} \varphi(\mathbf{x}_t)$ , and use it to perturb the predicted mean, steering the sampling step towards valid regions:

$$\mathbf{x}_{t-1} \sim \mathcal{N}\left(\mu_\theta(\mathbf{x}_t, t, \mathcal{F}) + \lambda \Sigma_\theta(\mathbf{x}_t, t, \mathcal{F}) \nabla_{\mathbf{x}_t} \varphi(\mathbf{x}_t, \mathcal{F}), \Sigma_\theta(\mathbf{x}_t, t, \mathcal{F})\right) \quad (2)$$

Iterating this process yields the final sample  $\mathbf{x}_0$ . The complete procedure is shown in Algorithm 1.

### 3.3. Object Quantity Constraint

To control object quantity, we introduce a differentiable guidance function,  $\varphi_{\text{quantity}}$ , which operates on the predicted class logits during the reverse diffusion process. Our scene representation uses  $N_{\text{max}}$  potential object slots, where each slot’s logits include a channel for an “empty” class, denoted  $c_i$ . To enforce a target count  $N_{\text{target}}$ , we define a binary target vector  $\mathbf{T} \in \mathbb{R}^{N_{\text{max}}}$  specifying which of the  $N_{\text{max}}$  slots should be non-empty.

#### Algorithm 1: Guidance Sampling in Model

**Modules:** Model  $p_\theta(\cdot|\mathcal{F})$ , guidance functions

$$\varphi(\cdot) = \{\varphi_{\text{quantity}}(\cdot), \varphi_{\text{articoll}}(\cdot)\}.$$

1 // constraint-guided learning

**Input:** 3D scene layout  $\mathbf{x} = \{\mathbf{o}_1, \dots, \mathbf{o}_N\}$  and floor plan  $\mathcal{F}$ , where  $N$  is a fixed number of objects.

2 **repeat**

3    $\mathbf{x}_0 \sim p(\mathbf{x}_0|\mathcal{F})$

4    $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), t \sim \mathcal{U}(\{1, \dots, T\})$

5    $\mathbf{x}_t = \sqrt{\hat{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \hat{\alpha}_t} \epsilon, \tilde{\mathbf{x}}_0^t \sim p_\theta(\cdot)$

6    $\theta = \theta - \eta \nabla_\theta \|\epsilon - \epsilon_\theta(\mathbf{x}_t, t) - \lambda \Sigma \mathbf{g}\|_2^2$

7 **until converged;**

8 // one-step guided sampling

9 **function sample** ( $\tau^t, \varphi$ ):

10    $\mu = \mu_\theta(\mathbf{x}_t, t, \mathcal{F}), \Sigma = \Sigma_\theta(\mathbf{x}_t, t, \mathcal{F})$

11    $\varphi(\mathbf{x}_t) = \gamma_1 \varphi_{\text{quantity}}(\mathbf{x}_t) + \gamma_2 \varphi_{\text{articoll}}(\mathbf{x}_t)$

12    $\mathbf{x}_{t-1} = \mathcal{N}(\mathbf{x}_{t-1}; \mu + \lambda \Sigma \nabla_{\mathbf{x}_t} \varphi(\mathbf{x}_t, \mathcal{F}) |_{\mathbf{x}_t=\mu}, \Sigma)$

13   **return**  $\mathbf{x}_{t-1}$

14 // constraint-guided generation

**Input:** initial scene layout  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

15 **for**  $t = T, \dots, 1$  **do**

16   // sampling with optimization

17    $\mathbf{x}_{t-1} = \text{sample}(\mathbf{x}_t, \varphi)$

18 **end**

19 **return**  $\mathbf{x}_0$

The guidance function is formulated as the Binary Cross-Entropy (BCEWithLogits) loss between the predicted “empty” logits and this target vector:

$$\varphi_{\text{quantity}}(\mathbf{x}) = \text{BCEWithLogits}(\{c_i\}_{i=1}^{N_{\text{max}}}, \mathbf{T}) \quad (3)$$

The gradient of this function,  $\nabla_{\mathbf{x}_t} \varphi_{\text{quantity}}$ , provides a direct signal during sampling, steering the model to populate the scene with the  $N_{\text{target}}$  desired objects.



### 3.4. Articulated Collision Constraint

Standard collision losses are insufficient as they only check static geometry, ignoring functional plausibility. We introduce a differentiable guidance function,  $\varphi_{\text{articoll}}$ , to penalize such "functional collisions." Our method quantifies functional collision during diffusion sampling. For each object  $b_i$  in the scene  $\mathcal{B}$ , we identify if it is articulated via a lookup. If so, we compute its *functionally extended state*,  $b'_i$ , by heuristically expanding its bounding box along its primary axis of articulation. For non-articulated objects,  $b'_i = b_i$ . The total collision penalty is the sum of pairwise 3D Intersection over Union (IoU) between each object's extended state  $b'_i$  and all other static objects  $b_j$ :

$$\varphi_{\text{articoll}}(\mathbf{x}) = \sum_{i=1}^N \sum_{j=1, j \neq i}^N \text{IoU}_{3D}(b'_i, b_j) \quad (4)$$

This penalty is differentiable with respect to scene parameters. Its gradient,  $\nabla_{\mathbf{x}_t} \varphi_{\text{articoll}}$ , steers the reverse diffusion process away from obstructed configurations, ensuring generated scenes are functionally viable.

### 3.5. Walkable Area Control

Ensuring a specific walkable space ratio is critical for robotic navigation. Enforcing such a constraint directly within the diffusion sampling loop would be computationally prohibitive, likely requiring expensive spatial queries at every step, and could potentially destabilize the generative process. We therefore introduce an efficient post-processing optimization as shown in Algorithm 2, which decouples semantic layout generation from spatial density tuning, as shown in Figure 4. Our algorithm iteratively refines the scene to meet a target ratio  $\tau$  by modifying only object *sizes* while preserving their *placements*. This strategy retains the core semantic structure while guaranteeing navigability.

### 3.6. Proposed Task-Specific Evaluation Metrics

To measure the controllability of Sec. 3.1, Sec. 3.3, Sec. 3.4, and Sec. 3.5, we proposed four novel evaluation metrics to validate them. Including LLM Controllability (Sec. 3.6.1), Object Quantity Controllability (Sec. 3.6.2), Articulated Object Collision Ratio (Sec. 3.6.3) and Walkable Area Controllability (Sec. 3.6.4).

#### 3.6.1. LLM-Guided Layout Metric

To evaluate the structural and semantic fidelity of a generated graph  $G_{\text{gen}} = (V_{\text{gen}}, E_{\text{gen}})$  against the ground-truth  $G_{\text{gt}} = (V_{\text{gt}}, E_{\text{gt}})$ , we measure node similarity, constraint satisfaction, and edge similarity.

**Node Similarity.** We compute a maximum cardinality matching  $M : V_{\text{gen}} \rightarrow V_{\text{gt}}$  constrained by node type ( $T(v_{\text{gen}}) = T(v_{\text{gt}})$ ). The score is the match size normal-

---

#### Algorithm 2: Walkable Area Optimization

---

**Input:** Batch size  $B$ , max iterations  $M$ , ratio threshold  $\tau$ , top- $k$  objects  $k$

```

1 for  $i \leftarrow 1$  to  $B$  do
2   Extract data for scene  $i$ ;
3    $iter \leftarrow 0$ ;
4    $r_i \leftarrow \text{CalculateWalkableRatio}(\text{scene } i)$ ;
5   while  $r_i < \tau$  and  $iter < M$  do
6      $iter \leftarrow iter + 1$ ;
7     Sort valid objects by area;
8      $replacement \leftarrow \text{False}$ ;
9     for top  $k$  objects do
10      Find closest object in database;
11      if smaller replacement found then
12        Replace size and features;
13         $replacement \leftarrow \text{True}$ ;
14      end
15    end
16  end
17  if not replacement then
18    break;
19  end
20   $r_i \leftarrow \text{CalculateWalkableRatio}(\text{scene } i)$ ;
21 end
22 return  $optimized\_scenes$ 

```

---

ized by the larger graph size to penalize extraneous nodes:

$$S_{\text{node}}(G_{\text{gen}}, G_{\text{gt}}) = \frac{|M|}{\max(|V_{\text{gen}}|, |V_{\text{gt}}|)} \quad (5)$$

**Constraint Satisfaction Score.** This metric evaluates the area ratio distribution per room type ( $R(G, c)$ ). We first measure the L1 distance between the generated and ground-truth distributions:

$$D_{L1} = \sum_{c \in \mathcal{C}} |R(G_{\text{gen}}, c) - R(G_{\text{gt}}, c)| \quad (6)$$

The normalized constraint satisfaction score  $S_{\text{constraint}} \in [0, 1]$  is defined as:

$$S_{\text{constraint}}(G_{\text{gen}}, G_{\text{gt}}) = 1 - \frac{1}{2} D_{L1} \quad (7)$$

**Edge Similarity.** Based on the node matching  $M$ , we identify the set of matched edges  $E_{\text{match}}$ , where edges in  $E_{\text{gen}}$  have corresponding nodes (under  $M$ ) that are also connected in  $E_{\text{gt}}$ :

$$E_{\text{match}} = \{(u, v) \in E_{\text{gen}} \mid (M(u), M(v)) \in E_{\text{gt}}\} \quad (8)$$

The score is normalized by the larger edge set to penalize spurious edges:

$$S_{\text{edge}}(G_{\text{gen}}, G_{\text{gt}}) = \frac{|E_{\text{match}}|}{\max(|E_{\text{gen}}|, |E_{\text{gt}}|)} \quad (9)$$

### 3.6.2. Object Quantity Control Metric

To assess the model’s capability to control the number of objects within a generated scene, we conducted a quantitative evaluation. We prompted the model to generate rooms containing a specific target number of objects,  $N_{\text{target}}$ . A scene  $S_i$  is considered a “success” if its generated object count, denoted  $N(S_i)$ , exactly matches the target. The Success Rate (SR) for a given target quantity over a set of  $M$  scenes is then formally defined as:

$$SR(N_{\text{target}}) = \frac{1}{M} \sum_{i=1}^M \mathbb{I}(N(S_i) = N_{\text{target}}) \quad (10)$$

where  $\mathbb{I}(\cdot)$  is the indicator function. For each target quantity, we calculated this rate over  $M = 100$  generated scenes. The results, presented in Table 3, show that our model maintains a high SR, consistently above 95% for most tested quantities, demonstrating precise control over scene composition.

### 3.6.3. Articulation Collision Metric

To quantitatively validate the effectiveness of our proposed articulated object collision constraint,  $\varphi_{\text{articoll}}$ , we conduct an evaluation by comparing our guided model against a baseline variant trained without this constraint. For each model, we generate a test set of 100 scenes and post-process them by transforming all articulated objects into their functionally extended states to simulate real-world usage.

We introduce  $R_{\text{acoll}}$  as our primary metric, defined as the proportion of articulated objects involved in functional collisions. Formally, for a given scene  $S$ , it is calculated as:

$$R_{\text{acoll}}(S) = \frac{1}{N_A} \sum_{j \in \mathcal{A}} \mathbb{I} \left( \max_{i \neq j} (\text{IoU}_{3D}(b'_i, b'_j)) > 0 \right) \quad (11)$$

where  $\mathcal{A}$  is the set of articulated objects in the scene ( $N_A = |\mathcal{A}|$ ), and  $b'_i, b'_j$  represent the functional bounding boxes. For an articulated object, this is its volume in the extended state; for a static object, it is its original bounding box. The indicator function  $\mathbb{I}(\cdot)$  returns 1 if the condition is true. A lower  $R_{\text{acoll}}$  score indicates superior functional plausibility.

### 3.6.4. Walkable Area Controllability Metric

To quantitatively measure the navigability and spaciousness of a generated scene, we define the  $R_{\text{walkable}}$ . This metric is calculated as the ratio of the total unobstructed floor area to the total area of the room. Let  $A_{\text{room}}$  be the total area

of the floor plan. For a scene containing  $N$  objects, where the floor footprint of the  $i$ -th object is denoted by  $A_i$ , the total walkable area  $A_{\text{walkable}}$  is the room area minus the sum of all object footprints. The ratio is formally defined as:

$$R_{\text{walkable}} = \frac{A_{\text{walkable}}}{A_{\text{room}}} = \frac{A_{\text{room}} - \sum_{i=1}^N A_i}{A_{\text{room}}} \quad (12)$$

The SR over a set of  $M$  generated scenes for a given threshold is then formally defined as:

$$SR(\tau_{\text{walkable}}) = \frac{1}{M} \sum_{i=1}^M \mathbb{I}(R_{\text{walkable}}(S_i) \geq \tau_{\text{walkable}}) \quad (13)$$

## 4. Experiment

### 4.1. Implementation Detail

**Datasets.** Our pipeline uses three specialized datasets. The layout generator is trained on 3D-FRONT [8] with 14,629 indoor scenes, and layouts are populated with textured assets from 3D-FUTURE [7] containing 16,563 furniture models. To enable interactivity, we use GPartNet [9], which provides part-level semantics and pose data for 8,489 parts across 1,166 objects.

**Baselines.** We compare with three baselines to demonstrate methodological progression. ATISS [16] is an autoregressive Transformer for sequential object generation, Dif-fuScene [19] employs diffusion to improve global consistency, and PhyScene [20] adds physics-based guidance for physically plausible scene synthesis.

**Evaluation Metrics.** We assess layout quality using Fréchet Inception Distance (FID) [10], Kernel Inception Distance (KID) [2], Scene Classification Accuracy (SCA), and Category KL divergence (CKL) following [19]. Constraint-specific metrics are defined in Sec. 3.6 and evaluated in Sec. 4.3–4.6.

### 4.2. Conditioned Scene Synthesis Evaluation

Conditioned scene synthesis is evaluated with diverse textual and spatial prompts. As shown in Figure 5, Scene-Foundry generates layouts that align with user-defined conditions and maintain spatial coherence. Quantitative results in Table 1 show high structural realism, minimal artifacts, and superior KID and CKL scores, confirming controllable high-fidelity scene generation.

### 4.3. LLM-Guided Layout Generation Evaluation

The resulting layouts were quantitatively evaluated against ground-truth graphs that perfectly satisfy the given high-level constraints. We use our proposed room graph similarity metrics ( $S_{\text{node}}$ ,  $S_{\text{constraint}}$ , and  $S_{\text{edge}}$ ) from Sec. 3.6.1 to measure the semantic and structural fidelity of the generated floor plans. Our method achieves extremely high scores across all three metrics as shown in Table 2.



Figure 5. Qualitative comparison of conditioned scene synthesis results among PhyScene, ATISS, DiffuScene, and SceneFoundry.

Table 1. **Floor-conditioned Scene Synthesis.** We compare SceneFoundry with baseline on common perceptual quality scores FID, KID, SCA, CKL.

Method	FID ↓	KID ↓	SCA	CKL ↓
ATISS	30.19	0.0010	49.14	0.0028
DiffuScene	<b>25.00</b>	<b>0.0004</b>	51.78	0.0031
PhyScene	25.52	0.0006	<b>50.10</b>	0.0025
SceneFoundry	29.02	<b>0.0004</b>	49.11	<b>0.0024</b>

Table 2. Results for our LLM Controllability experiment.

Method	$S_{\text{node}} \uparrow$	$S_{\text{constraint}} \uparrow$	$S_{\text{edge}} \uparrow$
Ours (LLM control)	<b>0.989</b>	<b>0.923</b>	<b>0.954</b>

#### 4.4. Object Quantity Controllability Evaluation

Rooms are generated with target object counts  $N_{\text{target}}$  ranging from 5 to 16. A generation is successful when the final count of non-empty slots matches  $N_{\text{target}}$ . For each target, 100 scenes are sampled, and the success rate (SR) is reported in Table 3. Results show consistently high SR values (0.95–0.97), demonstrating stable quantity control and robustness to scene complexity.

#### 4.5. Articulated Collision Constraint Evaluation

We further evaluate the effect of the proposed Articulated Object Collision Constraint, which enforces functional clearance for movable furniture. As illustrated in

Table 3. SR of generating scenes with a specific target number of objects. The rate is calculated over 100 generated scenes for each target.

$N_{\text{target}}$	SR	$N_{\text{target}}$	SR	$N_{\text{target}}$	SR
5	0.95	9	0.96	13	0.96
6	0.95	10	0.97	14	0.95
7	0.96	11	0.96	15	0.95
8	0.95	12	0.96	16	0.95

Table 4. Comparison of  $R_{\text{acoll}}$  and  $R_{\text{reach}}$ . Our method drastically reduces functional collisions ( $R_{\text{acoll}} \downarrow$ ) and improves object accessibility ( $R_{\text{reach}} \uparrow$ ).

Method	$R_{\text{acoll}} \downarrow$	$R_{\text{reach}} \uparrow$
Baseline (w/o $\varphi_{\text{articoll}}$ )	0.191	0.742
Ours (w/ $\varphi_{\text{articoll}}$ )	<b>0.109</b>	<b>0.808</b>

Figure 6, scenes generated without this constraint often contain obstructed articulated parts, such as drawers or chairs that cannot move freely. When the constraint is applied, these collisions are effectively eliminated, resulting in functionally usable layouts. Quantitatively, our method achieves a significantly lower functional collision rate  $R_{\text{acoll}}$  (Sec. 3.6.3) and higher object reachability ( $R_{\text{reach}}$ ) [20] than the baseline, as shown in Table 4. The constraint improves scene functionality and accessibility in generated scenes.



Figure 6. Visualization of the Articulated Object Collision Constraint. Synthesized scenes without the constraint (top) show obstructed articulated furniture, such as drawers that cannot open, while applying the constraint (bottom) enables proper motion and functional layouts.

#### 4.6. Walkable Area Controllability Evaluation

Thresholds from 0.60 to 0.95 were tested with  $M = 100$  scenes. For each threshold, we compare the SR under two conditions: a baseline without our constraint and with our constraint activated. Our method significantly increases the SR across all tested thresholds, as shown in Figure 7. Qualitative examples in Figure 8 further show that the constraint maintains sufficient free space for navigation while preserving realistic scene density.

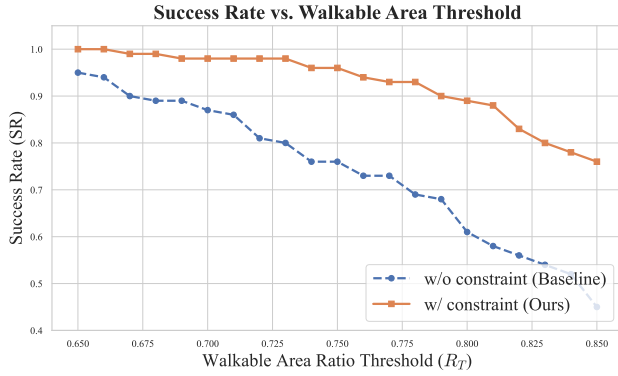


Figure 7. Success Rate (SR) versus Walkable Area Ratio Threshold ( $R_T$ ). Walkable Area Control (orange) consistently outperforms the baseline (blue), ensuring navigable layouts.

#### 4.7. Ablation Study on Scene Generation Control

We conduct an ablation study to validate the effectiveness of our proposed guidance mechanisms for controlling scene plausibility. Following the evaluation protocol established by PhyScene [20], we measure three key metrics: object collision ( $\text{Col}_{\text{obj}} \downarrow$ ), walkable ratio ( $\text{R}_{\text{walkable}} \uparrow$ ), and object reachability ( $\text{R}_{\text{reach}} \uparrow$ ). We analyze the contributions of our two guidance functions: ArtiCollision and Walkable Ratio, as shown in Table 5.

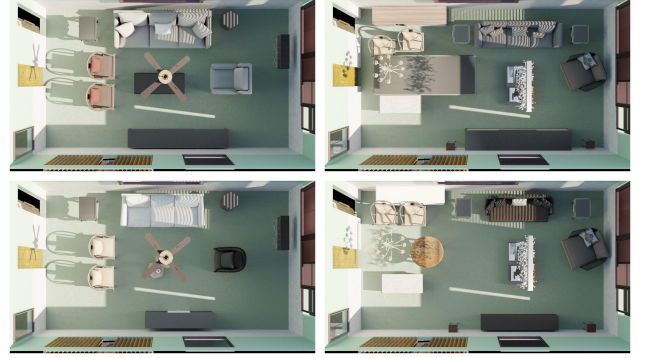


Figure 8. Visualization of Walkable Area Control.

Table 5. Ablation study on the use of guidance functions. Our final result balances the effectiveness of the two guidance mechanisms.

ArtiCollision	Walkable Ratio	$\text{Col}_{\text{obj}} \downarrow$	$\text{R}_{\text{walkable}} \uparrow$	$\text{R}_{\text{reach}} \uparrow$
✓		0.279	0.774	0.742
		0.267	0.774	0.808
	✓	0.250	0.822	0.782
✓	✓	<b>0.249</b>	<b>0.822</b>	<b>0.830</b>

## 5. Conclusion

This paper presents a multi-stage framework for controllable 3D environment generation that connects user intent with structured scene synthesis. The method integrates an LLM-guided floor-plan generator and a diffusion-based layout model to achieve semantic and spatial coherence. Functional modules including Object Quantity Control, Articulated Object Collision, and Walkable Area Control ensure realistic, accessible, and navigable layouts. Experimental results demonstrate precise control, high perceptual quality, and strong functional plausibility, consistently outperforming baseline methods and establishing a solid foundation for future Sim-to-Real research.



## References

- [1] Arpit Bansal, Hong-Min Chu, Avi Schwarzschild, Soumyadip Sengupta, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Universal guidance for diffusion models, 2023. [2](#)
- [2] Mikołaj Bińkowski, Danica J. Sutherland, Michael Arbel, and Arthur Gretton. Demystifying mmd gans, 2021. [6](#), [11](#)
- [3] Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Nießner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from rgb-d data in indoor environments, 2017. [2](#)
- [4] Zoey Chen, Sho Kiami, Abhishek Gupta, and Vikash Kumar. Genau: Retargeting behaviors to unseen situations via generative augmentation, 2023. [2](#)
- [5] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes, 2017. [2](#)
- [6] Prafulla Dhariwal and Alex Nichol. Diffusion models beat gans on image synthesis, 2021. [2](#)
- [7] Huan Fu, Rongfei Jia, Lin Gao, Mingming Gong, Binqiang Zhao, Steve Maybank, and Dacheng Tao. 3d-future: 3d furniture shape with texture, 2020. [6](#)
- [8] Huan Fu, Bowen Cai, Lin Gao, Lingxiao Zhang, Jiaming Wang Cao Li, Zengqi Xun, Chengyue Sun, Rongfei Jia, Binqiang Zhao, and Hao Zhang. 3d-front: 3d furnished rooms with layouts and semantics, 2021. [2](#), [6](#)
- [9] Haoran Geng, Helin Xu, Chengyang Zhao, Chao Xu, Li Yi, Siyuan Huang, and He Wang. Gapartnet: Cross-category domain-generalizable object perception and manipulation via generalizable and actionable parts, 2023. [2](#), [3](#), [6](#)
- [10] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium, 2018. [6](#), [11](#)
- [11] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance, 2022. [2](#)
- [12] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models, 2020. [2](#), [10](#), [11](#)
- [13] Siyuan Huang, Zan Wang, Puhao Li, Baoxiong Jia, Tengyu Liu, Yixin Zhu, Wei Liang, and Song-Chun Zhu. Diffusion-based generation, optimization, and planning in 3d scenes, 2023. [2](#)
- [14] Eric Kolve, Roozbeh Mottaghi, Winson Han, Eli VanderBilt, Luca Weihs, Alvaro Herrasti, Matt Deitke, Kiana Ehsani, Daniel Gordon, Yuke Zhu, Aniruddha Kembhavi, Abhinav Gupta, and Ali Farhadi. Ai2-thor: An interactive 3d environment for visual ai, 2022. [2](#)
- [15] Zixi Liang, Guowei Xu, Haifeng Wu, Ye Huang, Wen Li, and Lixin Duan. S-inf: Towards realistic indoor scene synthesis via scene implicit neural field, 2025. [2](#)
- [16] Despoina Paschalidou, Amlan Kar, Maria Shugrina, Karsten Kreis, Andreas Geiger, and Sanja Fidler. Atiss: Autoregressive transformers for indoor scene synthesis, 2021. [2](#), [6](#)
- [17] Alexander Raistrick, Lingjie Mei, Karhan Kayan, David Yan, Yiming Zuo, Beining Han, Hongyu Wen, Meenal Parakh, Stamatis Alexandropoulos, Lahav Lipson, Zeyu Ma, and Jia Deng. Infinigen indoors: Photorealistic indoor scenes using procedural generation, 2024. [2](#), [3](#)
- [18] Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, Devi Parikh, and Dhruv Batra. Habitat: A platform for embodied ai research, 2019. [2](#)
- [19] Jiapeng Tang, Yinyu Nie, Lev Markhasin, Angela Dai, Justus Thies, and Matthias Nießner. Diffuscene: Denoising diffusion models for generative indoor scene synthesis, 2024. [2](#), [3](#), [6](#), [11](#)
- [20] Yandan Yang, Baoxiong Jia, Peiyuan Zhi, and Siyuan Huang. Physcene: Physically interactable 3d scene synthesis for embodied ai, 2024. [2](#), [3](#), [6](#), [7](#), [8](#)

# Appendix

## SCENEFOUNDRY: Generating Interactive Infinite 3D Worlds

### A. Reproducibility and Code Release

To ensure the reproducibility of our results and facilitate future research in controllable scene generation, we will release our complete source code, trained model weights, and the post-processing scripts upon acceptance. Detailed instructions for environment setup and inference are provided in the supplementary material.

<https://github.com/anc891203/SceneFoundry>

### B. Discussion

We provide a critical analysis of the current limitations of our **SceneFoundry** framework and discuss the broader societal implications of our work below.

#### B.1. Limitation

Despite the demonstrated efficacy of SceneFoundry in generating functionally viable and apartment-scale environments, several limitations remain to be addressed.

**Inference Latency.** A primary constraint is the computational cost associated with the multi-stage pipeline. While the LLM-guided floor plan generation is relatively efficient, the core diffusion-based furniture population relies on iterative denoising steps. Coupled with the gradient calculations required for our novel constraints, the inference time for a full apartment scale is considerable. This currently precludes the system from real-time generation applications.

**Heuristic Approximation of Articulation.** Our Articulated Object Collision Constraint relies on a heuristic expansion of bounding boxes to approximate the kinematic workspace of objects. While robust for standard furniture morphologies found in GPartNet, this axis-aligned expansion simplifies the complex, potentially non-linear trajectories of certain articulated parts. Consequently, for highly complex mechanisms or multi-jointed objects, the collision avoidance might be overly conservative or, in rare cases, insufficient.

**Dataset Bias and Generalization.** The stylistic and semantic diversity of our generated scenes is inherently

bounded by the underlying training data, specifically 3D-FRONT and GPartNet. While these datasets are extensive, they may not fully encompass the architectural styles of different cultures or historical periods. As with all learning-based generative models, the system may exhibit biases present in the dataset, potentially favoring modern, Western-style interior layouts over others.

#### B.2. Social Impact

The primary societal contribution of this work lies in its potential to accelerate the development of embodied AI and service robotics. By automating the synthesis of large-scale, functionally sound training environments, SceneFoundry significantly reduces the reliance on costly and labor-intensive real-world data collection. This democratization of high-quality simulation data can foster innovation in domestic robotics, potentially leading to deploying intelligent agents capable of assisting the elderly or individuals with disabilities in their daily lives.

From an environmental perspective, while the training of large diffusion models incurs a carbon footprint, the ability to train robots in simulation (Sim-to-Real) drastically reduces the energy consumption, material waste, and physical risks associated with trial-and-error learning in the physical world. Regarding ethical considerations, unlike generative models for faces or media, the generation of indoor scene layouts carries a relatively low risk of malicious misuse. However, we advocate for continued awareness regarding the cultural biases embedded in synthetic datasets to ensure the inclusivity of future AI technologies.

### C. Preliminaries

Our generative framework is built upon Denoising Diffusion Probabilistic Models (DDPMs) [12]. DDPMs are a class of latent variable models designed to learn a data distribution  $p(\mathbf{x})$  by reversing a gradual noising process. In this section, we briefly review the mathematical formulation of DDPMs, including the forward diffusion process, the reverse denoising process, and the training objective.

#### C.1. Denoising Diffusion Probabilistic Models

**Forward Diffusion Process.** The forward process, also known as the diffusion process, is a fixed Markov chain

that gradually adds Gaussian noise to the data  $\mathbf{x}_0 \sim q(\mathbf{x}_0)$  over a sequence of timesteps  $t = 1, \dots, T$ . The transition probability at each step is defined as:

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}), \quad (14)$$

where  $\{\beta_t \in (0, 1)\}_{t=1}^T$  is a pre-defined variance schedule. As  $T \rightarrow \infty$ , the data  $\mathbf{x}_T$  approaches an isotropic Gaussian distribution  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ . A key property of this process is that we can sample  $\mathbf{x}_t$  at any arbitrary timestep  $t$  directly from  $\mathbf{x}_0$  in closed form:

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}), \quad (15)$$

where  $\alpha_t = 1 - \beta_t$  and  $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$ . This allows us to express  $\mathbf{x}_t$  as a linear combination of the original data and noise:

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}). \quad (16)$$

**Reverse Denoising Process.** The goal of the generative model is to reverse this process, sampling from  $q(\mathbf{x}_{t-1} | \mathbf{x}_t)$  to reconstruct the data. Since the exact posterior is intractable, we approximate it using a learnable Markov chain with parameterized Gaussian transitions:

$$p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t)), \quad (17)$$

starting from  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . The mean  $\boldsymbol{\mu}_\theta$  and covariance  $\boldsymbol{\Sigma}_\theta$  are predicted by neural networks. Following [12], we set  $\boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t) = \sigma_t^2 \mathbf{I}$ , where  $\sigma_t^2$  is set to  $\beta_t$  or  $\tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t$ . The mean is parameterized to predict the noise  $\epsilon$  added to  $\mathbf{x}_0$ :

$$\boldsymbol{\mu}_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right). \quad (18)$$

**Optimization Objective.** The model is trained by optimizing the variational lower bound on the negative log-likelihood. Ho et al. [12] demonstrated that a simplified objective yields better sample quality. This simplified loss calculates the mean squared error between the true noise  $\epsilon$  and the predicted noise  $\epsilon_\theta$ :

$$\mathcal{L}_{\text{simple}} = \mathbb{E}_{t, \mathbf{x}_0, \epsilon} [\|\epsilon - \epsilon_\theta(\mathbf{x}_t, t)\|^2], \quad (19)$$

where  $t$  is uniformly sampled from  $\{1, \dots, T\}$ . In our framework, we adapt this backbone to generate 3D scene layouts by conditioning the denoiser on floor plan constraints.

## D. Implementation Details

**Experimental Setting.** We evaluate our method on the 3D-FRONT dataset, employing the official train/test splits to ensure consistency with prior work. For articulation-aware generation, we augment the object assets using the GPartNet dataset, which provides part-level annotations. To verify the robustness of our method, all baselines are retrained on this identical data subset. We generate 1,000 scenes for each experimental condition to compute reliable metrics.

**Evaluation Metrics.** To quantitatively evaluate the quality, diversity, and semantic coherence of our generated scenes, we employ a comprehensive suite of metrics. *Standard Perceptual & Semantic Metrics:*

- **Fréchet Inception Distance (FID) [10]:** Measures the distributional distance between deep features of generated ( $\mu_g, \Sigma_g$ ) and real ( $\mu_r, \Sigma_r$ ) scene renderings:

$$\text{FID} = \|\mu_r - \mu_g\|^2 + \text{Tr}(\Sigma_r + \Sigma_g - 2(\Sigma_r \Sigma_g)^{1/2}) \quad (20)$$

- **Kernel Inception Distance (KID) [2]:** An unbiased estimator of the Maximum Mean Discrepancy (MMD) between feature representations, suitable for smaller sample sizes:

$$\begin{aligned} \text{KID} &= \text{MMD}^2(P_r, P_g) \\ &= \mathbb{E}[k(x, x')] + \mathbb{E}[k(y, y')] - 2\mathbb{E}[k(x, y)] \end{aligned} \quad (21)$$

- **Scene-Class Alignment (SCA) [19]:** Evaluates semantic consistency by calculating the classification accuracy of a pre-trained scene classifier  $C$  on generated layouts  $x_{\text{gen}}$ :

$$\text{SCA} = \mathbb{E}_{x_{\text{gen}}} [\mathbb{I}(C(x_{\text{gen}}) = y_{\text{label}})] \quad (22)$$

- **Category KL divergence (CKL) [19]:** Measures the divergence between the object category distribution of the generated set ( $P_g$ ) and the ground truth ( $P_r$ ):

$$\text{CKL} = D_{KL}(P_g || P_r) = \sum_i P_g(i) \log \frac{P_g(i)}{P_r(i)} \quad (23)$$

*Proposed Controllability Metrics (Ours):*

- **LLM-Guided Layout Metric:** Evaluates the structural and semantic fidelity of the generated floor plan graph against the ground truth by assessing node matching, edge connectivity, and constraint satisfaction.
- **Object Quantity Control Metric:** Defines the success rate of generating scenes that contain the exact target number of objects specified by the user.

- **Articulation Collision Ratio:** Measures the percentage of articulated objects (e.g., cabinets) that are functionally obstructed by other objects when in their open/extended state.
- **Walkable Area Controllability:** Calculates the success rate of generating scenes where the ratio of unobstructed walkable floor area meets or exceeds a specified threshold.

### D.1. Compare Model Settings

We benchmark SceneFoundry against three state-of-the-art baselines:

- **ATISS:** An autoregressive transformer model that places objects sequentially. We use the official implementation, retraining it on our dataset split for fair comparison.
- **DiffuScene:** A diffusion-based model that generates scene layouts in parallel. This represents the current state-of-the-art in unconditional layout generation.
- **PhyScene:** A recent physics-aware generative model. We compare against PhyScene to highlight the advantages of our specific articulated object constraints.

All models receive the same floor plan input during the conditional generation tasks.

### D.2. Training Details

The model is trained using the Adam optimizer with a learning rate of  $2 \times 10^{-4}$  and weight decay of 0.0. We employ a step learning rate schedule with a step size of 20,000 and a decay factor of 0.5. Training runs for 130,000 epochs with a batch size of 128. The gradient norm is clipped at 10. Table 6 summarizes the complete hyperparameter configuration.

Table 6. Detailed hyperparameter settings for training the diffusion model.

Configuration	Value
Optimizer	Adam
Base Learning Rate	$2 \times 10^{-4}$
Weight Decay	0.0
Batch Size	128
Max Gradient Norm	10
Learning Rate Schedule	Step Decay
LR Step Size	20,000
LR Decay Factor ( $\gamma$ )	0.5
Total Epochs	130,000

### D.3. Computing Resource Configuration

All model training and evaluation were conducted on a computing node equipped with a single **NVIDIA 3090 GPU (24GB VRAM)** and an **Intel Core i9-12900K**. Under this configuration, training the core diffusion model takes approximately 1500 hours. During inference, generating a complete apartment-scale scene (3 rooms) with full constraint guidance takes approximately 300 seconds per scene.

## E. Additional Experiments

### E.1. LLM Controllability

To rigorously evaluate the fidelity of our LLM-based parameter space guidance, we designed a comprehensive benchmark consisting of 20 distinct natural language prompts.

**Overall Performance.** We define a generation as successful only if the final layout strictly adheres to all constraints specified in the prompt, as shown in Table 7.

Table 7. Summary of LLM Controllability Experiments.

Metric	Value
Total Test Cases	20
Average Score	96.5%

**Component Analysis.** To understand the specific challenges in layout control, we decompose the performance into three sub-metrics. The score is composed of weights in the Table 8.

- **Room Presence:** Existence of required room types.
- **Adjacency:** Correct connectivity between rooms.
- **Constraints:** Geometric or functional requirements.

Table 8. Weights indicate the relative importance assigned to each component in the overall score.

Component	Avg. Score	Weight
Room Presence	98.9%	50%
Adjacency	92.3%	30%
Constraints	95.4%	20%

**Detailed Test Results** We provide a granular breakdown of each test case in Table 9. In these cases, the system maintains high scores on connectivity and functional constraints, ensuring the generated layouts remain usable.



Table 9. Complete Test Results for LLM Controllability. This table details the performance of 20 distinct test prompts.

Test ID	Score	Time(s)	Room Presence	Adjacency	Constraints
two_bedroom_aprt.01	100.0%	262.2	100%	100%	100%
open_plan_loft.01	100.0%	274.7	100%	100%	100%
family_home.01	100.0%	276.7	100%	100%	100%
master_suite_home.01	100.0%	265.3	100%	100%	100%
four_bedroom_house.01	100.0%	270.8	100%	100%	100%
guest_suite_home.01	100.0%	282.2	100%	100%	100%
dual_master_suite.01	100.0%	269.5	100%	100%	100%
balcony_apartment.01	100.0%	285.2	100%	100%	100%
multigenerational_home.01	100.0%	289.7	100%	100%	100%
basic_studio.01	100.0%	259.0	100%	100%	100%
one_bedroom_aprt.01	100.0%	270.0	100%	100%	100%
compact_efficiency.01	100.0%	269.7	100%	100%	100%
student_apartment.01	100.0%	266.0	100%	100%	100%
single_floor_accessible.01	100.0%	262.5	100%	100%	100%
entertainment_home.01	95.8%	274.9	100%	100%	79.2%
work_from_home.01	90.0%	263.0	80.0%	100%	100%
luxury_penthouse.01	88.8%	270.8	87.5%	100%	75.0%
three_bed_townhouse.01	85.0%	268.3	100%	50%	100%
separated_zones.01	85.0%	285.7	100%	50.0%	100%
home_office_layout.01	62.5%	267.8	25.0%	100%	100%

## F. Render Results

We present an extensive qualitative evaluation of the 3D indoor layouts generated by our proposed method. While quantitative metrics provide numerical evidence of our model’s performance, visual inspection is equally crucial for assessing the perceptual quality, spatial coherence, and practical usability of the synthesized scenes. To this end, we provide a comprehensive gallery of results across a diverse range of scene categories, demonstrating the robustness of our approach in handling complex room geometries.

To guarantee the highest visual fidelity, all visualizations were produced using the **Blender** creation suite, leveraging its advanced **Cycles** rendering engine. Cycles is a production-grade, physically-based path tracer that excels at simulating the intricate interactions of light transport. Unlike real-time rasterization engines, Cycles calculates global illumination, multi-bounce indirect lighting, and accurate soft shadows, which are essential for verifying that objects are properly grounded and not floating. Furthermore, we utilized high-resolution Physically Based Rendering textures and materials to enhance the realism of

the furniture, allowing for a rigorous assessment of the layout’s aesthetic quality.

The rendering pipeline imposes significant computational demands, particularly when processing scenes with high-poly assets and complex lighting setups. Consequently, all rendering tasks were executed on a dedicated workstation equipped with an **Intel Core i3-14100F CPU** and an **NVIDIA RTX 5090 GPU**. The massive 32GB VRAM of the RTX 5090 proved instrumental in loading large-scale scene data and high-resolution textures without memory bottlenecks, while the CPU efficiently managed scene graph traversal and asset loading.

These visualizations highlight the model’s capability to handle complex spatial arrangements naturally. We observe that the generated objects are physically plausible, exhibiting proper orientations and avoiding inter-object collisions, as shown in Figure 9 and Figure 10. Collectively, these qualitative results validate that our method not only adheres to rigid geometric constraints but also produces aesthetically pleasing and functionally realistic environments suitable for practical design applications.



Figure 9. **Generated 3D Layouts.** Representative visualization of scenes generated by SceneFoundry. (Part 1).



Figure 10. **Generated 3D Layouts.** Representative visualization of scenes generated by SceneFoundry. (Part 2).