

# Diffusion-based Generation, Optimization, and Planning in 3D Scenes

Siyuan Huang<sup>1\*</sup>, Zan Wang<sup>1,2\*</sup>, Puhao Li<sup>1,3</sup>, Baoxiong Jia<sup>1</sup>  
Tengyu Liu<sup>1</sup>, Yixin Zhu<sup>4</sup>, Wei Liang<sup>2</sup>, Song-Chun Zhu<sup>1,3,4</sup>

<sup>1</sup> National Key Laboratory of General Artificial Intelligence, BIGAI

<sup>2</sup> School of Computer Science & Technology, Beijing Institute of Technology

<sup>3</sup> Dept. of Automation, Tsinghua University    <sup>4</sup> Institute for AI, Peking University

<https://scenediffuser.github.io>



Figure 1. **Illustration of the SceneDiffuser**, applicable to various scene-conditioned 3D tasks: (a) **human pose generation**, (b) **human motion generation**, (c) **dexterous grasp generation**, (d) **path planning for 3D navigation with goals**, and (e) **motion planning for robot arms**.

## Abstract

We introduce *SceneDiffuser*, a conditional generative model for 3D scene understanding. *SceneDiffuser* provides a unified model for solving scene-conditioned generation, optimization, and planning. In contrast to prior work, *SceneDiffuser* is intrinsically scene-aware, physics-based, and goal-oriented. With an iterative sampling strategy, *SceneDiffuser* jointly formulates the scene-aware generation, physics-based optimization, and goal-oriented planning via a diffusion-based denoising process in a fully differentiable fashion. Such a design alleviates the discrepancies among different modules and the posterior collapse

of previous scene-conditioned generative models. We evaluate *SceneDiffuser* on various 3D scene understanding tasks, including human pose and motion generation, dexterous grasp generation, path planning for 3D navigation, and motion planning for robot arms. The results show significant improvements compared with previous models, demonstrating the tremendous potential of *SceneDiffuser* for the broad community of 3D scene understanding.

## 1. Introduction

The ability to generate, optimize, and plan in 3D scenes is a long-standing goal for multiple research domains across computer vision, graphics, and robotics. Various tasks have been devised to achieve these goals, fostering downstream applications in motion generation [32, 67, 70, 83], motion planning [40, 41, 58, 72], grasp generation [25, 31, 33], nav-

\* These authors contributed equally to this work.

 Corresponding authors: Siyuan Huang (syhuang@bigai.ai) and Wei Liang (liangwei@bit.edu.cn).

igation [1, 88], embodied perception and manipulation [30, 37, 59], and autonomous driving [3, 50].

Despite rich applications and great successes, existing models designed for these tasks exhibit two fundamental limitations for real-world 3D scene understanding.

First, most prior work [8, 14, 32, 47, 58, 66–68, 70] leverages the conditional Variational Autoencoder (cVAE) for the conditional generation in 3D scenes. cVAE model utilizes an encoder-decoder structure to learn the posterior distribution and relies on the learned latent variables to sample. Although cVAE is easy to train and sample due to its simple architecture and one-step sampling procedure, it suffers from the **posterior collapse problem** [12, 17, 26, 62, 67, 77, 86]; the learned latent variable is ignored by a strong decoder, leading to limited generation diversity from these collapsed modes. Such collapse is further magnified in 3D tasks with stronger 3D decoders and more complex and noisy input conditions, *e.g.*, the natural 3D scans [9].

Second, despite the close relations among generation, optimization, and planning in 3D scenes, there **lacks a unified framework** that could address existing discrepancies among these models. Previous work [15, 33, 67] applies off-the-shelf physics-based post-optimization methods over outputs of generative models and often produces inconsistent and implausible generations, especially when transferring to novel scenes. Similarly, planners are usually standalone modules over results of generative model [8, 14] for trajectory planning or learned separately with the reinforcement learning (RL) [88], leading to gaps between planning and other modules (*e.g.*, generation) during inference, especially in novel scenes where explorations are limited.

To tackle these limitations, we introduce SceneDiffuser, a conditional generative model based on the diffusion process. SceneDiffuser eliminates the discrepancies and provides a single home for scene-conditioned generation, optimization, and planning. Specifically, with a denoising process, it learns a diffusion model for scene-conditioned generation while training. During inference, SceneDiffuser jointly solves the scene-aware generation, physics-based optimization, and goal-oriented planning through a unified iterative guided-sampling framework. Such a design equips SceneDiffuser with the following superiority:

**Generation** Building upon the diffusion model, SceneDiffuser significantly alleviates the posterior collapse problem of scene-conditioned generative models. Since the forward diffusion process can be treated as data augmentation in 3D scenes, it helps traverse sufficient scene-conditioned distribution modes.

**Optimization** SceneDiffuser integrates the physics-based objective into each step of the sampling process as conditional guidance, enabling the differentiable physics-based optimization during both the learning and sampling

process. This design facilitates the physically-plausible generation, which is critical for tasks in 3D scenes.

**Planning** Based on the scene-conditioned trajectory-level generator, SceneDiffuser possesses a global trajectory planner with physics and goal awareness, making the learned planner generalize better to long-horizon trajectories and novel 3D scenes.

As illustrated in Fig. 1, we evaluate SceneDiffuser on diverse 3D scene understanding tasks. The results on human pose, motion, and dexterous grasp generation significantly improve, demonstrating plausible and diverse generations with 3D scene and object conditions. The results on path planning for 3D navigation and motion planning for robot arms reveal the generalizable and long-horizon planning capability of SceneDiffuser.

## 2. Related Work

**Conditional Generation in 3D Scenes** Generating diverse contents and rich interactions in 3D scenes is essential for understanding the 3D scene affordances. Recently, we have witnessed several applications on conditional scene generation [24, 45, 69, 81], human pose [16, 32, 80, 83, 85] and motion generation [14, 32, 47, 58, 66–68, 70] in furnished 3D indoor scenes, and object-conditioned grasp pose generation [25, 31, 33, 60, 74]. However, most previous methods [6, 14, 16, 25, 31, 60, 66, 70, 72] rely on cVAE and suffer from the posterior collapse problem [12, 17, 26, 62, 67, 77, 86], especially when the 3D scene is natural and complex. In this work, SceneDiffuser addresses the posterior collapse with the diffusion-based denoising process.

**Physics-based Optimization in 3D Scenes** Producing physically plausible generations compatible with 3D scenes is one of the challenges in the scene-conditioned generation. Previous work uses physics-based post-optimization [15, 33, 67] or differentiable objective [25, 70, 83] to integrate collision and contact constraints into the generation framework. However, post-optimization approaches [15, 33, 67] are oftentimes inefficient and cannot be learned jointly with the generative models, yielding inconsistent generation results. Similarly, differentiable approaches [25, 70, 83] post constraints on the final objective, thus cannot optimize the physical interactions during the sampling, producing implausible generations, especially when adapting to novel scenes. In this work, SceneDiffuser eliminates such inconsistency with the differentiable physics-based optimization integrated into each step of the sampling process.

**Planning in 3D Scenes** The ability to act and plan in 3D scenes is critical for an intelligent agent and has led to the recent culmination of embodied AI research [28, 30, 52, 53]. Among all tasks, visual navigation has been most studied in the vision and robotics community [4, 13, 21, 38, 73, 87, 88]. However, existing works rely heavily on model-based planning with the single-step dy-

dynamic model [5, 11, 46, 65, 71, 75], lacking a trajectory-level optimization for long-horizon planning. Further, the physical interactions are not explicitly modeled into the planning. This deficiency makes it challenging to generalize to natural scenes, where exploration is limited, and fast learning and adaptation are required. In comparison, with the global trajectory planner based on a trajectory-level generator, SceneDiffuser demonstrates better generalization in long-horizon plans and novel 3D scenes.

**Diffusion-based Models** Diffusion model [19, 22, 55, 57] has come forth into a promising class of generative model for learning and sampling data distributions with an iterative denoising process, facilitating the image [10, 56], text [76], and shape generation [34]. With flexible conditioning, it is further extended to the language-conditioned image [39, 49, 51], video [18, 54], and 3D generation [42, 61, 78]. Notably, Janner *et al.* [23] integrate the generation and planning into the same sampling framework for behavior synthesis. To our best knowledge, SceneDiffuser is the first framework that models the 3D scene-conditioned generation with a diffusion model and integrates the generation, optimization, and planning into a unified framework.

### 3. Background

#### 3.1. Problem Definition

Given a 3D scene  $\mathcal{S}$ , we aim to generate the optimal solution for completing the tasks (*e.g.*, navigation, manipulation) given the goal  $\mathcal{G}$  in the scene. We denote the state and action of an agent as  $(\mathbf{s}, \mathbf{a})$ . The dynamic model defines the state transition as  $p(\mathbf{s}_{i+1}|\mathbf{s}_i, \mathbf{a}_i)$ , which is often deterministic in scene understanding (*i.e.*,  $f(\mathbf{s}_i, \mathbf{a}_i)$ ). The trajectory is defined as  $\tau = (\mathbf{s}_0, \mathbf{a}_0, \dots, \mathbf{s}_i, \mathbf{a}_i, \dots, \mathbf{s}_N)$ , where  $N$  denotes the horizon of task solving in discrete time.

#### 3.2. Planning with Trajectory Optimization

The scene-conditional trajectory optimization is defined as maximizing the task objective:

$$\tau^* = \arg \max_{\tau} \mathcal{J}(\tau|\mathcal{S}, \mathcal{G}). \quad (1)$$

The dynamic model is usually known for trajectory optimization. Considering the future actions and states with predictable dynamics, the entire trajectory  $\tau$  can be optimized jointly and non-progressively with traditional [29] or data-driven [7] planning algorithms. Trajectory-based optimization benefits from its global awareness of history and future states, thus can better model the long-horizon tasks compared with single-step models in RL, where  $\mathbf{a}_{0:N}^* = \arg \max_{\mathbf{a}_{0:N}} \sum_{i=0}^N r(\mathbf{s}_i, \mathbf{a}_i|\mathcal{S}, \mathcal{G})$ .

#### 3.3. Diffusion Model

Diffusion models [19, 22, 55] are a class of generative models that represent the data generation with an it-

erative denoising process from Gaussian noise. It consists of a forward and a reverse process. The forward process  $q(\tau^t|\tau^{t-1})$  gradually destroys data  $\tau^0 \sim q(\tau^0)$  into Gaussian noise. The parametrized reverse process  $p_{\theta}(\tau^{t-1}|\tau^t)$  recovers the data from noise with the learned normal distribution from a fixed timestep. The training objective for  $\theta$  is denoising score matching over multiple noise scale [22, 64]. Please refer to the Appendix A for detailed descriptions of the diffusion model and its variants.

### 4. SceneDiffuser

SceneDiffuser models planning as trajectory optimization and solves the aforementioned problem with the spirit of *planning as sampling*, where the trajectory optimization is achieved by sampling trajectory-level distribution learned by the model. Leveraging the diffusion model with gradient-based sampling and flexible conditioning, SceneDiffuser models the scene-conditioned goal-oriented trajectory  $p(\tau^0|\mathcal{S}, \mathcal{G})$ :

$$p(\tau^0|\mathcal{S}, \mathcal{G}) = \frac{p_{\theta}(\tau^0|\mathcal{S})p_{\phi}(\mathcal{G}|\tau^0, \mathcal{S})}{p(\mathcal{G}|\mathcal{S})} \propto p_{\theta}(\tau^0|\mathcal{S})p_{\phi}(\mathcal{G}|\tau^0, \mathcal{S}). \quad (2)$$

**Generation**  $p_{\theta}(\tau^0|\mathcal{S})$  characterizes the probability of generating certain trajectories with the scene condition. It can be modeled using a conditional diffusion model [19, 55] with an iterative denoising process:

$$p_{\theta}(\tau^0|\mathcal{S}) = p(\tau^T) \prod_{t=1}^T p(\tau^{t-1}|\tau^t, \mathcal{S}), \quad (3)$$

$$p(\tau^{t-1}|\tau^t, \mathcal{S}) = \mathcal{N}(\tau^{t-1}; \mu_{\theta}(\tau^t, t, \mathcal{S}), \Sigma_{\theta}(\tau^t, t, \mathcal{S})).$$

**Optimization and Planning**  $p_{\phi}(\mathcal{G}|\tau^0, \mathcal{S})$  represents the probability of reaching the goal with the sampled trajectory, where the goal can be flexibly defined by customized objective functions in various tasks. As shown in Eq. (4), the precise definition of this probability is  $p_{\phi}(\mathcal{O} = 1|\tau^0, \mathcal{S}, \mathcal{G})$ , where  $\mathcal{O}$  is an optimality indicator that represents if the goal were achieved. Intuitively, the trajectory objective in Eq. (1) can be a good indicator for such optimality. We therefore expand  $p_{\phi}(\mathcal{G}|\tau^t, \mathcal{S})$  as its exponential in Eq. (5):

$$p_{\phi}(\mathcal{G}|\tau^t, \mathcal{S}) = p_{\phi}(\mathcal{O} = 1|\tau^t, \mathcal{S}, \mathcal{G}) \quad (4)$$

$$\propto \exp(\mathcal{J}(\tau^t|\mathcal{S}, \mathcal{G})) \quad (5)$$

$$= \exp(\varphi_p(\tau^t|\mathcal{S}, \mathcal{G}) + \varphi_o(\tau^t|\mathcal{S})). \quad (6)$$

Here,  $\varphi_o(\tau^t|\mathcal{S})$  denotes the objective for optimizing the trajectory with scene condition and is independent of task goal  $\mathcal{G}$ . In scene understanding,  $\varphi_o$  usually denotes plausible physical relationships (*e.g.*, collision, contact, and intersection).  $\varphi_p(\tau^t|\mathcal{S}, \mathcal{G})$  indicates the objective for planning (*i.e.*, goal-reaching) with scene condition. Both  $\varphi_o$  and  $\varphi_p$  can be explicitly defined or implicitly learned from observed trajectories with proper parametrization.

## 4.1. Learning

$p_\theta(\tau^0|\mathcal{S})$  is the scene-conditioned generator, which can be learned by the conditional diffusion model with the simplified objective of estimating the noise  $\epsilon$  [10, 19, 20], where

$$\begin{aligned} \mathcal{L}_\theta(\tau^0|\mathcal{S}) &= \mathbb{E}_{t, \epsilon, \tau^0} \left[ \|\epsilon - \epsilon_\theta(\sqrt{\hat{\alpha}^t} \tau^0 + \sqrt{1 - \hat{\alpha}^t} \epsilon, t, \mathcal{S})\|^2 \right] \\ &= \mathbb{E}_{t, \epsilon, \tau^0} \left[ \|\epsilon - \epsilon_\theta(\tau^t, t, \mathcal{S})\|^2 \right], \end{aligned} \quad (7)$$

where  $\hat{\alpha}^t$  is the pre-determined function in the forward process. With the learned  $p_\theta(\tau^0|\mathcal{S})$ , we sample  $p(\tau^0|\mathcal{S}, \mathcal{G})$  by taking the advantage of the diffusion model’s flexible conditioning [10, 23]. Specifically, we approximate  $p_\phi(\mathcal{G}|\tau^t, \mathcal{S})$  using the Taylor expansion around  $\tau^t = \mu$  at timestep  $t$  as

$$\log p_\phi(\mathcal{G}|\tau^t, \mathcal{S}) \approx (\tau^t - \mu) \mathbf{g} + C,$$

where  $C$  is a constant,  $\mu = \mu_\theta(\tau^t, t, \mathcal{S})$  and  $\Sigma = \Sigma_\theta(\tau^t, t, \mathcal{S})$  are the inferred parameters of original diffusion process, and

$$\begin{aligned} \mathbf{g} &= \nabla_{\tau^t} \log p_\phi(\mathcal{G}|\tau^t, \mathcal{S})|_{\tau^t = \mu} \\ &= \nabla_{\tau^t} (\varphi_o(\tau^t|\mathcal{S}) + \varphi_p(\tau^t|\mathcal{S}, \mathcal{G}))|_{\tau^t = \mu}. \end{aligned} \quad (8)$$

Therefore, we have

$$p(\tau^{t-1}|\tau^t, \mathcal{S}, \mathcal{G}) = \mathcal{N}(\tau^{t-1}; \mu + \lambda \Sigma \mathbf{g}, \Sigma), \quad (9)$$

where  $\lambda$  is the scaling factor for the guidance. With Eq. (9), we can sample  $\tau^t$  with the guidance of optimizing and planning objectives.

Of note,  $\varphi_p$  and  $\varphi_o$  serve as the pre-defined guidance for tilting the original trajectory with physical and goal constraints. However, they can also be learned from the observed trajectories. During training, we first fix the learned base model of  $p_\theta(\tau^0|\mathcal{S})$ , then learn  $\phi_o$  and  $\phi_p$  for optimization and planning with the following objective:

$$\mathcal{L}_\phi(\tau^0|\mathcal{S}, \mathcal{G}) = \mathbb{E}_{t, \epsilon, \tau^0} \left[ \|\epsilon - \epsilon_\theta(\tau^t, t, \mathcal{S}) - \Sigma \mathbf{g}\|^2 \right]. \quad (10)$$

Alg. 1 summarizes the training procedure.

## 4.2. Sampling

With different sampling strategies, SceneDiffuser can generate, optimize, and plan the trajectory in 3D scenes, under a unified framework of guided sampling. Alg. 2 summarizes the detailed sampling algorithm.

**Scene-aware Generation** Sampling  $\tau^0$  from the distribution  $p_\theta(\tau^0|\mathcal{S})$  in Eq. (3) directly solves the conditional generation tasks. The sampled trajectories represent diverse modes and possible interactions with the 3D scenes.

**Physics-based Optimization** The physical relations between each state and the environment are defined by  $\varphi_o$  in Eq. (4) in a differentiable manner. For general optimization without the planning objective, the task goal  $\mathcal{G}$  is to sample a plausible trajectory in 3D scenes. Therefore, we can draw physically plausible trajectories in 3D scenes by sampling from  $p(\tau^0|\mathcal{S}, \mathcal{G})$  with Eq. (9).

---

### Algorithm 1: Training of the SceneDiffuser

---

```

1 // train base generation model
  Input: Trajectory in 3D scene ( $\tau^0, \mathcal{S}$ )
2 repeat
3    $\tau^0 \sim p(\tau^0|\mathcal{S})$ 
4    $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), t \sim \mathcal{U}(\{1, \dots, T\})$ 
5    $\tau^t = \sqrt{\hat{\alpha}^t} \tau^0 + \sqrt{1 - \hat{\alpha}^t} \epsilon$ 
6    $\theta = \theta - \eta \nabla_\theta \|\epsilon - \epsilon_\theta(\tau^t, t, \mathcal{S})\|_2^2$ 
7 until converged;
8 // (optional) train optimization and
  planning model
  Input: Trajectory in 3D scene with goal ( $\tau^0, \mathcal{S}, \mathcal{G}$ ), learned  $\theta$ 
    for  $p_\theta(\tau^0|\mathcal{S})$ 
9 repeat
10   $\tau^0 \sim p(\tau^0|\mathcal{S}, \mathcal{G})$ 
11   $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), t \sim \mathcal{U}(\{1, \dots, T\})$ 
12   $\mu = \mu_\theta(\tau^t, t, \mathcal{S}), \Sigma = \Sigma_\theta(\tau^t, t, \mathcal{S})$ 
13   $\mathbf{g} = \nabla_{\tau^t} \log p_\phi(\mathcal{G}|\tau^t, \mathcal{S})|_{\tau^t = \mu}$ 
14   $\tau^t = \sqrt{\hat{\alpha}^t} \tau^0 + \sqrt{1 - \hat{\alpha}^t} \epsilon$ 
15   $\phi = \phi - \eta \nabla_\phi \|\epsilon - \epsilon_\theta(\tau^t, t, \mathcal{S}) - \lambda \Sigma \mathbf{g}\|_2^2$ 
16 until converged;

```

---



---

### Algorithm 2: Sampling of the SceneDiffuser for generation, optimization, and planning

---

```

Modules: Model  $p_\theta(\cdot|\mathcal{S})$ , optimization objective  $\varphi_o(\cdot|\mathcal{S})$ ,
  and planner objective  $\varphi_p(\cdot|\mathcal{S}, \mathcal{G})$ 
1 // one-step guided sampling
2 function sample ( $\tau^t, \mathcal{J}$ ):
3    $\mu = \mu_\theta(\tau^t, t, \mathcal{S}), \Sigma = \Sigma_\theta(\tau^t, t, \mathcal{S})$ 
4    $\tau^{t-1} = \mathcal{N}(\tau^{t-1}; \mu + \lambda \Sigma \nabla_{\tau^t} (\mathcal{J}(\tau^t|\mathcal{S}, \mathcal{G}))|_{\tau^t = \mu}, \Sigma)$ 
5   return  $\tau^{t-1}$ 
6 // physics-based generation
  Input: initial trajectory  $\tau^T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
7 for  $t = T, \dots, 1$  do
8   // sampling with optimization
9    $\tau^{t-1} = \text{sample}(\tau^t, \varphi_o(\cdot|\mathcal{S}))$ 
10 return  $\tau^0$ 
11 // goal-oriented planning
  Input: planning steps  $N$ , starting state  $\hat{\mathbf{s}}_0$ , initial plan
     $\tau_0^T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
12  $i = 1$ 
13 while not done and planning step  $i < N$  do
14   for  $t = T, \dots, 1$  do
15      $\tau_i^{t-1} = \text{sample}(\tau_i^t, \varphi_o(\cdot|\mathcal{S}) + \varphi_p(\cdot|\mathcal{S}, \mathcal{G}))$ 
16     // planning as inpainting
17      $\tau_i^{t-1}[0:i] = \hat{\mathbf{s}}_{0:i}$ 
18   Act  $\hat{a}_{i-1}$  to reach  $\hat{\mathbf{s}}_i = \tau_i^0[i], \hat{\mathbf{s}}_{0:i} = \hat{\mathbf{s}}_{0:i-1} \cup \hat{\mathbf{s}}_i$ 
19   Increment planning step  $i = i + 1$ 

```

---

**Goal-oriented Planning** The goal-oriented planning can be formulated as motion inpainting under the sampling framework. Given the start state  $\hat{\mathbf{s}}_s$  and the goal state  $\hat{\mathbf{s}}_g$ , the planning module returns trajectory  $\hat{\tau} = (\hat{\mathbf{s}}_0, \hat{\mathbf{a}}_0, \dots, \hat{\mathbf{s}}_i, \hat{\mathbf{a}}_i, \dots, \hat{\mathbf{s}}_g)$  that can reach the goal state. We set the first state as  $\hat{\mathbf{s}}_0 = \hat{\mathbf{s}}_s$  and define the goal state and reward of goal-reaching in  $\varphi_p$ . For each step  $i$ , we first keep the previous states and inpaint the remaining trajectory by sampling the goal-oriented SceneDiffuser with an iterative denoising process. Next, we take the action that

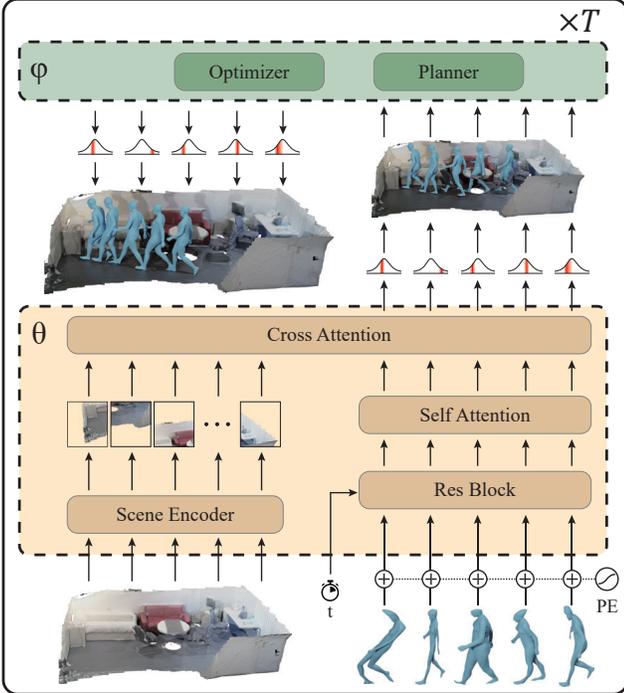


Figure 2. **Model architecture of the SceneDiffuser.** We use cross-attention to learn the relation between the input trajectory and scene condition. The optimizer and planner serve as the guidance for physically-plausible and goal-oriented trajectories.

can reach the next sampled state with  $(\hat{\mathbf{a}}_{i-1}, \hat{\mathbf{s}}_i)$ . As illustrated in Alg. 2, we repeat the planning steps until reaching the goal or the maximal planning step. Our planner leverages the trajectory-level generator, thus more generalizable to long-horizon trajectories and novel scenes.

### 4.3. Model Architecture

The design of SceneDiffuser follows the practices of conditional diffusion model [20, 49, 51]. Specifically, we augment the time-conditional diffusion model with cross-attention [63] for flexible conditioning. As shown in Fig. 2, for each sampling step, the model computes the cross-attention between the 3D scene condition and input trajectory, wherein the key and value are learned from the condition, and the query is learned from the input trajectory. The computed vector is fed into a feed-forward layer to estimate the noise  $\epsilon$ . The 3D scene is processed by a scene encoder (*i.e.*, Point Transformer [84] or PointNet [43]). Please refer to the Appendix B for details.

### 4.4. Objective Design

For optimization and planning objectives discussed in Sec. 4, we consider two types of trajectory objectives: (i) trajectory-level objective, and (ii) the accumulation of step-wise objective. For optimization, we consider step-wise collision and contact objective, as well as trajectory level

smoothness objective, *i.e.*,  $\{\varphi_o^{\text{collision}}, \varphi_o^{\text{contact}}, \varphi_o^{\text{smoothness}}\}$ . For planning, we consider the accumulation of simple step-wise distance *i.e.*,  $\varphi_p^{L_2}$ . Please refer to the Appendix C for implementation details of our objective design. Empirically, we observe that parameterizing the objectives with timestep  $t$  and increasing the guidance during the last several diffusion steps will enhance the effect of guidance.

## 5. Experiments

To demonstrate SceneDiffuser is general and applicable to various scenarios, we evaluate SceneDiffuser on five scene understanding tasks. For generation, we evaluate the scene-conditioned human pose and motion generation and object-conditioned dexterous grasp generation. For planning, we evaluate the path planning for 3D navigation and motion planning for robot arms. We first introduce the compared methods used in our experiments, followed by detailed settings, results analyses, and ablative studies for each task. Due to the page limit, we refer to the Appendix for more details about the implementation, experimental settings, and additional results and ablations.

### 5.1. Compared Methods

For conditional generation tasks, we primarily compare SceneDiffuser with the widely-adopted cVAE model [25, 31, 32, 70, 83] and its variants. We also compare with strategies for optimizing the physics of the trajectory in the cVAE, including integrating into training as loss and plugging upon as the post-optimization. For planning, we compare with a stochastic planner learned by imitation learning using Behavior Cloning (BC) and a simple heuristic-based deterministic planner guided by  $L_2$  distance.

### 5.2. Human Pose Generation in 3D Scenes

**Setup** Scene-conditioned human pose generation aims to generate semantically plausible and physically feasible single-frame human bodies within the given 3D scenes. We evaluate the task on the 12 indoor scenes provided by PROX [15] and the refined version of PROX’S per-frame SMPL-X parameters from LEMO [79]. The input is the colored point cloud extracted by randomly downsampling the scene meshes provided in PROX. Training/testing splits are created following the literature [67, 83], resulting in  $\sim 53k$  frames in 8 scenes for training and others for testing.

**Metrics** We evaluate the physical plausibility of generated poses with both direct human evaluations and indirect collision and contact scores. For the direct measure, we randomly selected 1000 frames in the four test scenes and instructed seven participants to decide whether the generated human pose was plausible. We compute the mean percentage of plausible generation and term this metric as the plausible rate. For indirect measures, we report (i) the non-collision score of the generated human bodies by calculat-

Table 1. **Quantitative results of human pose generation in 3D scenes.** We report metrics for physical plausibility and diversity.

model	plausible rate $\uparrow$	non-collision score $\uparrow$	contact score $\uparrow$	APD (trans.) $\uparrow$	std (trans.) $\uparrow$	APD (param) $\uparrow$	std (param) $\uparrow$	APD (marker) $\uparrow$	std (marker) $\uparrow$
cVAE (w/o. $L_{HS}$ ) [83]	12.57	99.78	96.42	<b>1.218</b>	<b>0.494</b>	2.878	0.166	3.638	0.172
cVAE (w/ $L_{HS}$ ) [83]	14.64	99.75	99.25	1.013	0.416	2.994	0.170	3.614	0.169
our (w/o opt.)	24.83	99.74	<b>99.43</b>	0.776	0.331	3.204	0.195	3.483	0.167
our (w/ opt.)	<b>49.35</b>	<b>99.93</b>	98.05	1.009	0.413	<b>3.297</b>	<b>0.197</b>	<b>3.679</b>	<b>0.177</b>

Table 2. **Quantitative results of human motion generation in 3D scenes.** We report model variants with and without the start pose.

model	plausible rate $\uparrow$	non-collision score $\uparrow$	contact score $\uparrow$	APD (trans.) $\uparrow$	std (trans.) $\uparrow$	APD (param) $\uparrow$	std (param) $\uparrow$	APD (marker) $\uparrow$	std (marker) $\uparrow$
cVAE (w/o start) [70]	5.88	99.86	86.26	<b>1.628</b>	<b>0.613</b>	<b>2.766</b>	<b>0.155</b>	3.275	0.150
ours (w/o start)	<b>24.70</b>	<b>99.71</b>	<b>97.92</b>	0.568	0.237	2.339	0.126	3.299	0.151
ours (w/o start & w/ opt.)	23.53	99.70	97.84	0.542	0.226	2.338	0.125	<b>3.301</b>	<b>0.151</b>
cVAE (w/ start) [70]	16.24	<b>99.88</b>	95.44	<b>0.478</b>	<b>0.188</b>	<b>1.747</b>	<b>0.091</b>	<b>2.308</b>	<b>0.105</b>
Ours (w/ start)	41.76	99.85	<b>99.63</b>	0.193	0.081	1.372	0.065	1.568	0.072
Ours (w/ start & w/ opt.)	<b>42.30</b>	99.85	99.62	0.192	0.080	1.368	0.063	1.565	0.075

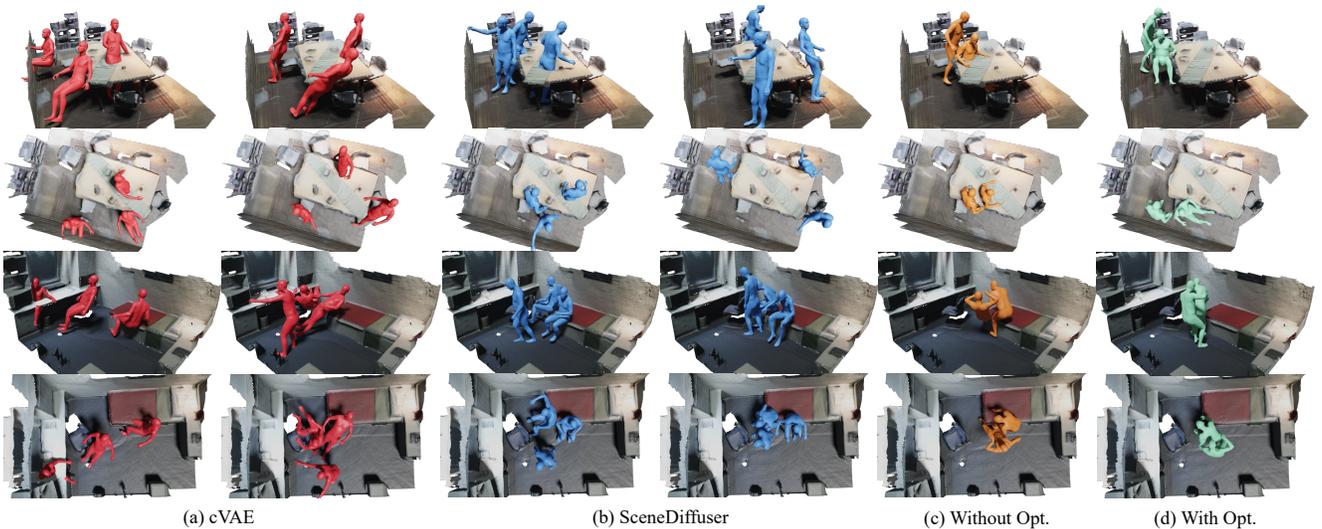


Figure 3. **Qualitative results of human pose generation in 3D scenes.** From left to right: (a) cVAE generation, (b) SceneDiffuser generation without optimization, and poses generated (c) with and (d) without applying our optimization-guided sampling.

ing the proportion of the scene vertices with positive SDF to the human body and (ii) the contact score by checking if the body contact with the scene in a distance [15] below a pre-defined threshold. Following the literature [77, 83], we evaluate the diversities of global translation, generated SMPL-X parameters, and the marker-based body-mesh representation [82]. Specifically, we calculate the diversity of generated pose with the Average Pairwise Distance (APD) and standard deviation (std).

**Results** Tab. 1 quantitatively demonstrates that SceneDiffuser generates significantly better poses while maintaining generation diversity. We further provide qualitative comparisons between baseline models and SceneDiffuser in Fig. 3. While achieving a comparable performance of diversity, collision, and contact, our model generates results that contain considerably more physically plausible poses (*e.g.*, floating, severe collision). This is reflected by the significant superiority (*i.e.*, over 30%) over

cVAE-based baselines on plausible rates. We observe this large improvement both quantitatively from the plausible rate and non-collision score and qualitatively in Fig. 3. Notably, our optimization-guided sampling improves the generator with 25% on the plausible rate, showing the efficacy of the proposed optimization-guided sampling strategy and its potential for a broader range of 3D tasks with physic-based constraints or objectives.

### 5.3. Human Motion Generation in 3D Scenes

**Setup** We consider generating human motion sequences under two different settings: (1) condition solely on the 3D scene, and (2) condition on both the starting pose and the 3D scene. We use the same human and scene representation as in Sec. 5.2 and clip the original LEMO motion sequence into segments with a fixed duration (60 frames). In total, we obtain 28k motion segments with the distance between each start and end pose being longer than 0.2 me-



Figure 4. **Human motions generated by SceneDiffuser.** Each row shows sampled human motions from the same start pose.

ters. We follow the same split in Sec. 5.2 for training/testing and the same evaluation metrics for the pose generation. We report the average values of pose metrics over motion sequence as our performance measure.

**Results** As quantitatively shown in Tab. 2, SceneDiffuser consistently generates high-quality motion sequences compared to cVAE baselines. Specifically, our generated motion outperforms baseline models on both plausible rate and contact scores. This performance gain indicates better coverage of motion that involves rich interaction with the scene while remaining physically plausible. It also causes lower diversity in metrics (*e.g.*, translation variance) since the plausible space for the motion is limited compared with cVAE. Empirically, we observe that providing the start position of motion as a condition constrains possible future motion sequences and leads to a drop in generation diversity for all models. In addition, providing the start condition benefits the physical plausibility since the motion starts from a plausible pose. We also note only a marginal performance improvement after applying optimization-guided sampling. One potential reason is that the generated motions are already plausible and receive small guidance from the optimization. As qualitatively shown in Fig. 4, SceneDiffuser generates diverse motions (*e.g.*, “sit,” “walk”) from the same start position in unseen 3D scenes.

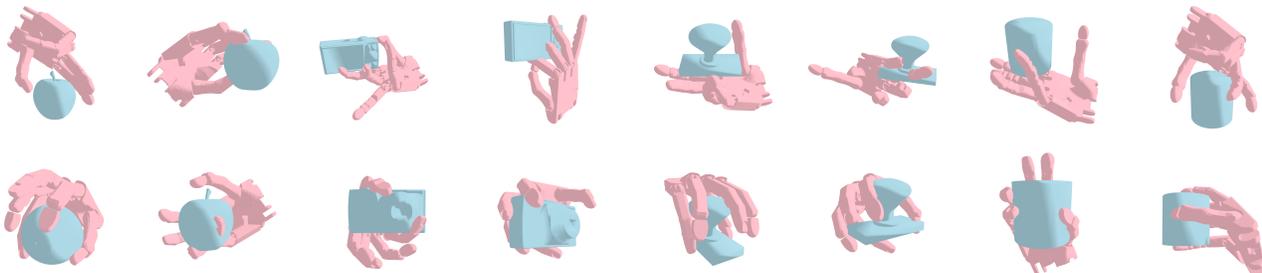


Figure 5. **Qualitative results of dexterous grasp generation.** Compared to grasps generated by cVAE (first row), SceneDiffuser (second row) generates fewer colliding or floating poses, which helps to achieve a higher success rate.

Table 3. **Quantitative results of dexterous grasp generation on MultiDex [31] dataset.** We measure the success rates under different diversities and depth collisions. TTA. denotes test-time optimization with physics and contact.

model	succ. rate (%) $\uparrow$			depth coll. (mm) $\downarrow$
	$\sigma$	$2\sigma$	all	
cVAE [25]	0.00	10.09	14.06	22.98
cVAE (w/ TTA.) [25]	0.00	21.91	17.97	15.19
ours (w/o opt.)	70.65	<b>71.25</b>	<b>71.25</b>	17.34
ours (w/ opt.)	<b>71.27</b>	69.84	69.84	<b>14.61</b>

#### 5.4. Dexterous Grasp Generation for 3D Objects

**Setup** Dexterous grasp generation aims to generate diverse and stable grasping poses for the given object with a human-like dexterous hand. We use the Shadowhand subset of the MultiDex [31] dataset, which contains diverse dexterous grasping poses for 58 daily objects. We represent the pose of Shadowhand as  $q := (t, R, \theta) \in \mathbb{R}^{33}$ , where  $t \in \mathbb{R}^3$  and  $R \in \mathbb{R}^6$  denote the global translation and orientation respectively, and  $\theta \in \mathbb{R}^{24}$  describes the rotation angles of the revolute joints. An object is represented by its point cloud  $\mathcal{O} \in \mathbb{R}^{2048 \times 3}$ . We split the dataset into 48 seen objects and 10 unseen objects for training and testing, respectively.

**Metrics** We evaluate models in terms of success rate, diversity, and collision depth. We test if a grasp is successful in IsaacGym [35] by applying external forces to the object and measuring the movement of the object. To measure how learned models capture the diversity of successful grasping pose in the training data, we report the success rate of generated poses that lies at different variance levels from the mean pose of training data. We measure the collision depth as the maximum depth that the hand penetrates the object in each successful grasp for testing models’ performance on physically correct grasps. In all cases, we ignore the root transformation of the hand as it does not contribute to the diversity of grasping types.

**Results** Tab. 3 quantitatively demonstrates that SceneDiffuser generates significantly better grasp poses in terms of success rate while correctly balancing the diversity

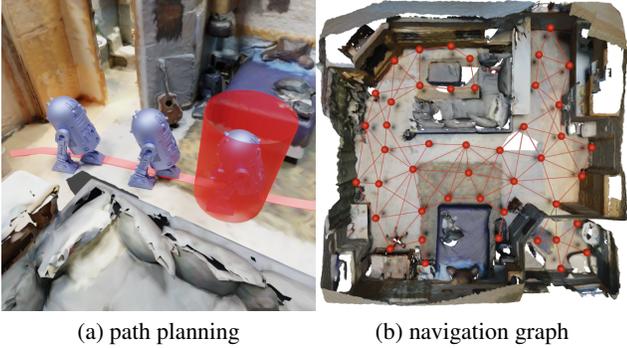


Figure 6. **Path planning for 3D scene navigation.** SceneDiffuser generates trajectories in long-horizon tasks.

of generation and grasp success. This result indicates that the SceneDiffuser achieves a consistently high success rate without much performance drop when the generated pose diverges from the mean pose in the training data. We also show that, by applying optimizer upon SceneDiffuser, the guided sampling process can reduce the violation of physically implausible grasping poses, outperforming the state-of-the-art baseline [25] without additional training or intermediate representation (*i.e.*, contact maps). We provide qualitative results in Fig. 5 for visualization.

### 5.5. Path Planning for 3D Scene Navigation

**Setup** We manually selected 61 indoor scenes from ScanNet [9] to construct room-level scenarios for navigational path planning and annotated these scenes with navigation graphs. As shown in Fig. 6b, these annotations are more spatially dense and physically plausible compared to previous methods [1]. We represent the physical robot with a cylinder to simulate physically plausible trajectories; see Fig. 6a. In total, we collected around 6k trajectories by searching the shortest paths between the randomly selected start and target nodes on the graph. We use trajectories in 46 scenes for training and trajectories in the rest 15 scenes for evaluation. Models take the input as the scene point cloud  $\mathcal{S} \in \mathbb{R}^{32768 \times 3}$ , a given start position  $\hat{s}_0 \in \mathbb{R}^2$ , and a target position  $\mathcal{G} \in \mathbb{R}^2$  on the floor plane.

**Metrics** We evaluate the planned results by checking if the “robot” can move from the start to the target without collision along the planned trajectory. We report the average success rate and planning steps over all test cases.

**Results** As shown in Tab. 4, SceneDiffuser outperforms both the BC and the deterministic planner baseline. These results indicate the efficacy of guided sampling with the planning objective, especially given that all test scenes are unseen during training. Crucially, as simple heuristics (like  $L_2$ ) oftentimes lead to dead-ends in path planning, SceneDiffuser can correctly combine past knowledge on the scene-conditioned trajectory distribution and planning ob-

Table 4. **Quantitative results of path planning in 3D navigation and motion planning for robot arms.**

task	model	succ. rate(%) $\uparrow$	planning steps $\downarrow$
path plan	BC	0	150
	deterministic( $L_2$ )	13.50	137.98
	ours	<b>73.75</b>	<b>90.38</b>
arm motion	BC	0.31	299.08
	deterministic( $L_2$ )	72.87	<b>141.28</b>
	ours	<b>78.59</b>	147.60

jective under specific unseen scenes to redirect planning direction, which helps to avoid obstacles and dead-ends to reach the goal successfully. Compared with the baseline models, our model also requires fewer planning steps while maintaining a higher success rate. This suggests that SceneDiffuser successfully navigates to the target without diverging even in long-horizon tasks, where classic RL-based stochastic planners suffer (*i.e.*, the low performance of BC).

### 5.6. Motion Planning for Robot Arms

**Setup** Aiming to generate valid robot arm motion trajectories in cluttered scenes, we used the Franka Emika arm with seven revolute joints and collected 19,800 trajectories over 200 randomly generated cluttered scenes using the MoveIt 2.0 [48], as shown in Fig. 7. We represent the scene with point clouds  $\mathcal{S} \in \mathbb{R}^{4096 \times 3}$  and the robot arm trajectory with a sequence of joint angles  $\mathcal{R} \in [-\pi, \pi]$ . We train our model on 160 scenes and test on 40 unseen scenes.

**Metrics** Similar to Sec. 5.5, we evaluate the generated trajectories by success rate on unseen scenes and the average number of planning steps. We consider a trajectory successful if the robot arm reaches the goal pose by a certain distance threshold within a limited number of steps.

**Results** We observe similar overall performance as in Sec. 5.5. Tab. 4 shows that SceneDiffuser consistently outperforms both the RL-based BC baseline and the deterministic planner baseline. SceneDiffuser’s planning steps for successful trials are also comparable with the deterministic planner, showing the efficacy of the planner in long-horizon scenarios.

### 5.7. Ablation Analyses

We explore how the scaling coefficient  $\lambda$  influences the human pose generation results. We report the diversity and physics metrics of sampling results under different  $\lambda$ s, ranging from 0.1 to 100. As shown in Tab. 5,  $\lambda$  balances generation collision/contact and diversity in human pose generation. Specifically,  $\lambda = 1.0$  leads to the best physical plausibility while larger  $\lambda$  values lead to diverse generation results. We attribute this effect to the optimization as with bigger  $\lambda$ s; the optimization will draw poses away from the

Table 5. Ablation of the scale coefficient for optimization.

metric	$\lambda = 0.1$	$\lambda = 1.0$	$\lambda = 10.0$	$\lambda = 100.0$
plausible rate $\uparrow$	28.75	<b>52.5</b>	21.25	0
APD (trans.) $\uparrow$	0.764	0.886	1.564	<b>23.96</b>
APD (param) $\uparrow$	3.206	3.243	9.040	<b>573.6</b>
non-collision score $\uparrow$	99.76	<b>99.87</b>	99.85	74.9
contact score $\uparrow$	<b>99.70</b>	99.65	81.75	0.0

scene. Due to the page limit, we provide more ablative studies in Appendix E, including the sampling steps, choices and hyperparameters of objectives, and model architectures.

### 5.8. Limitation

The primary limitation of the SceneDiffuser is its slow training and test speed compared to previous scene-conditioned generative models, a common issue of diffusion-based methods. We also observe that the optimization and planning are highly dependent on the objective designs, which requires efforts on hyper-parameter tuning.

## 6. Conclusion

We propose the SceneDiffuser as a general conditional generative model for generation, optimization, and planning in 3D scenes. SceneDiffuser is designed with appealing properties including scene-aware, physics-based, and goal-oriented. We demonstrate that the SceneDiffuser outperforms previous models by a large margin on various tasks, establishing its efficacy and flexibility.

A promising future direction is extending SceneDiffuser to richer 3D representations, including RGB-D images, semantic images, bird-eye view (BEV) images, videos, 3D meshes, and neural radiance field (NeRF) [36]. Such flexible conditions consume a tremendous amount of 3D training data, which is also a significant challenge. We also hope to extend the SceneDiffuser to outdoor scenes, e.g., the autonomous driving scenarios [3]. Moreover, the SceneDiffuser can be combined with recent large language

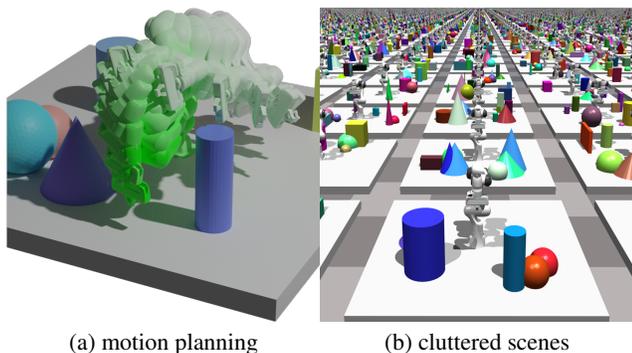


Figure 7. Motion planning for robot arms. SceneDiffuser generates arm motions in tabletop scenes with obstacles.

models (LLMs) [2] for automatic generation and planning with natural language instructions in 3D scenes, which is promising for the vision and robotics community. Finally, SceneDiffuser can serve as the tool for analyzing the behaviors of humans and agents if we can properly learn the planning objective, which naturally encodes the values and preferences that underlie the trajectories.

## 7. Acknowledgement

We thank Ruiqi Gao and Ying Nian Wu for their helpful discussions and suggestions. This work is supported in part by the National Key R&D Program of China (2021ZD0150200) and the Beijing Nova Program.

## References

- [1] Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian Reid, Stephen Gould, and Anton Van Den Hengel. Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2, 8
- [2] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 9
- [3] Sergio Casas, Abbas Sadat, and Raquel Urtasun. Mp3: A unified model to map, perceive, predict and plan. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2, 9
- [4] Devendra Singh Chaplot, Dhiraj Gandhi, Saurabh Gupta, Abhinav Gupta, and Ruslan Salakhutdinov. Learning to explore using active neural slam. In *International Conference on Learning Representations (ICLR)*, 2020. 2
- [5] Devendra Singh Chaplot, Dhiraj Prakashchand Gandhi, Abhinav Gupta, and Russ R Salakhutdinov. Object goal navigation using goal-oriented semantic exploration. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 3
- [6] Yuxiao Chen, Boris Ivanovic, and Marco Pavone. Scept: Scene-consistent, policy-based trajectory predictions for planning. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [7] Sanjiban Choudhury, Mohak Bhardwaj, Sankalp Arora, Ashish Kapoor, Gireeja Ranade, Sebastian Scherer, and Debadepta Dey. Data-driven planning via imitation learning. *International Journal of Robotics Research (IJRR)*, 37(13-14):1632–1672, 2018. 3
- [8] Alexander Cui, Sergio Casas, Abbas Sadat, Renjie Liao, and Raquel Urtasun. Lookout: Diverse multi-future prediction and planning for self-driving. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [9] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In

- Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2, 8, A3
- [10] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. 3, 4
- [11] Heming Du, Xin Yu, and Liang Zheng. Learning object relation graph and tentative policy for visual navigation. In *European Conference on Computer Vision (ECCV)*, 2020. 3
- [12] Hao Fu, Chunyuan Li, Xiaodong Liu, Jianfeng Gao, Asli Celikyilmaz, and Lawrence Carin. Cyclical annealing schedule: A simple approach to mitigating kl vanishing. In *North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, 2019. 2
- [13] Weituo Hao, Chunyuan Li, Xiujun Li, Lawrence Carin, and Jianfeng Gao. Towards learning a generic agent for vision-and-language navigation via pre-training. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2
- [14] Mohamed Hassan, Duygu Ceylan, Ruben Villegas, Jun Saito, Jimei Yang, Yi Zhou, and Michael J Black. Stochastic scene-aware motion prediction. In *International Conference on Computer Vision (ICCV)*, 2021. 2
- [15] Mohamed Hassan, Vasileios Choutas, Dimitrios Tzionas, and Michael J Black. Resolving 3d human pose ambiguities with 3d scene constraints. In *International Conference on Computer Vision (ICCV)*, 2019. 2, 5, 6, A1
- [16] Mohamed Hassan, Partha Ghosh, Joachim Tesch, Dimitrios Tzionas, and Michael J Black. Populating 3d scenes by learning human-scene interaction. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [17] Junxian He, Daniel Spokoyny, Graham Neubig, and Taylor Berg-Kirkpatrick. Lagging inference networks and posterior collapse in variational autoencoders. In *International Conference on Learning Representations (ICLR)*, 2018. 2
- [18] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022. 3
- [19] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 3, 4, A1, A3
- [20] Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded diffusion models for high fidelity image generation. *Journal of Machine Learning Research (JMLR)*, 23(47):1–33, 2022. 4, 5
- [21] Yicong Hong, Qi Wu, Yuankai Qi, Cristian Rodriguez-Opazo, and Stephen Gould. Vln bert: A recurrent vision-and-language bert for navigation. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [22] Aapo Hyvärinen and Peter Dayan. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research (JMLR)*, 6(4), 2005. 3
- [23] Michael Janner, Yilun Du, Joshua B Tenenbaum, and Sergey Levine. Planning with diffusion for flexible behavior synthesis. *arXiv preprint arXiv:2205.09991*, 2022. 3, 4
- [24] Chenfanfu Jiang, Siyuan Qi, Yixin Zhu, Siyuan Huang, Jenny Lin, Lap-Fai Yu, Demetri Terzopoulos, and Song-Chun Zhu. Configurable 3d scene synthesis and 2d image rendering with per-pixel ground truth using stochastic grammars. *International Journal of Computer Vision (IJCV)*, pages 920–941, 2018. 2
- [25] Hanwen Jiang, Shaowei Liu, Jiashun Wang, and Xiaolong Wang. Hand-object contact consistency reasoning for human grasps generation. In *International Conference on Computer Vision (ICCV)*, 2021. 1, 2, 5, 7, 8
- [26] Yoon Kim, Sam Wiseman, Andrew Miller, David Sontag, and Alexander Rush. Semi-amortized variational autoencoders. In *International Conference on Machine Learning (ICML)*, 2018. 2
- [27] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. A2
- [28] Eric Kolve, Roozbeh Mottaghi, Winson Han, Eli VanderBilt, Luca Weihs, Alvaro Herrasti, Daniel Gordon, Yuke Zhu, Abhinav Gupta, and Ali Farhadi. Ai2-thor: An interactive 3d environment for visual ai. *arXiv preprint arXiv:1712.05474*, 2017. 2
- [29] Steven M LaValle. *Planning algorithms*. Cambridge university press, 2006. 3
- [30] Chengshu Li, Fei Xia, Roberto Martín-Martín, Michael Lingelbach, Sanjana Srivastava, Bokui Shen, Kent Vainio, Cem Gokmen, Gokul Dharan, Tanish Jain, et al. igibson 2.0: Object-centric simulation for robot learning of everyday household tasks. In *Conference on Robot Learning (CoRL)*, 2021. 2
- [31] Puhao Li, Tengyu Liu, Yuyang Li, Yiran Geng, Yixin Zhu, Yaodong Yang, and Siyuan Huang. Gendexgrasp: Generalizable dexterous grasping. *arXiv preprint arXiv:2210.00722*, 2022. 1, 2, 5, 7, A2
- [32] Xueting Li, Sifei Liu, Kihwan Kim, Xiaolong Wang, Ming-Hsuan Yang, and Jan Kautz. Putting humans in a scene: Learning affordance in 3d indoor environments. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 1, 2, 5
- [33] Tengyu Liu, Zeyu Liu, Ziyuan Jiao, Yixin Zhu, and Song-Chun Zhu. Synthesizing diverse and physically stable grasps with arbitrary hand structures using differentiable force closure estimator. *IEEE Robotics and Automation Letters (RA-L)*, 7(1):470–477, 2021. 1, 2
- [34] Shitong Luo and Wei Hu. Diffusion probabilistic models for 3d point cloud generation. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 3
- [35] Viktor Makoviychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles Macklin, David Hoeller, Nikita Rudin, Arthur Allshire, Ankur Handa, et al. Isaac gym: High performance gpu based physics simulation for robot learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. 7, A3
- [36] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021.

- [37] Tongzhou Mu, Zhan Ling, Fanbo Xiang, Derek Yang, Xuanlin Li, Stone Tao, Zhiao Huang, Zhiwei Jia, and Hao Su. Maniskill: Generalizable manipulation skill benchmark with large-scale demonstrations. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. 2
- [38] Medhini Narasimhan, Erik Wijmans, Xinlei Chen, Trevor Darrell, Dhruv Batra, Devi Parikh, and Amanpreet Singh. Seeing the un-scene: Learning amodal semantic maps for room navigation. In *European Conference on Computer Vision (ECCV)*, 2020. 2
- [39] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021. 3, A1
- [40] Xue Bin Peng, Glen Berseth, and Michiel Van de Panne. Terrain-adaptive locomotion skills using deep reinforcement learning. *ACM Transactions on Graphics (TOG)*, 35(4):1–12, 2016. 1
- [41] Xue Bin Peng, Ze Ma, Pieter Abbeel, Sergey Levine, and Angjoo Kanazawa. Amp: Adversarial motion priors for stylized physics-based character control. *ACM Transactions on Graphics (TOG)*, 40(4):1–20, 2021. 1
- [42] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022. 3
- [43] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 5, A1, A3
- [44] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. A3
- [45] Siyuan Qi, Yixin Zhu, Siyuan Huang, Chenfanfu Jiang, and Song-Chun Zhu. Human-centric indoor scene synthesis using stochastic grammar. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2
- [46] Santhosh K Ramakrishnan, Dinesh Jayaraman, and Kristen Grauman. An exploration of embodied visual exploration. *International Journal of Computer Vision (IJCV)*, 129(5):1616–1649, 2021. 3
- [47] Davis Rempe, Tolga Birdal, Aaron Hertzmann, Jimei Yang, Srinath Sridhar, and Leonidas J Guibas. Humor: 3d human motion model for robust pose estimation. In *International Conference on Computer Vision (ICCV)*, 2021. 2
- [48] PickNik Robotics. Moveit motion planning framework for ros 2. <https://github.com/ros-planning/moveit2>, 2020. 8, A3
- [49] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 3, 5
- [50] Abbas Sadat, Mengye Ren, Andrei Pokrovsky, Yen-Chen Lin, Ersin Yumer, and Raquel Urtasun. Jointly learnable behavior and trajectory planning for self-driving vehicles. In *International Conference on Intelligent Robots and Systems (IROS)*, 2019. 2
- [51] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022. 3, 5
- [52] Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, et al. Habitat: A platform for embodied ai research. In *International Conference on Computer Vision (ICCV)*, 2019. 2
- [53] Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi, Luke Zettlemoyer, and Dieter Fox. Alfred: A benchmark for interpreting grounded instructions for everyday tasks. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2
- [54] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, et al. Make-a-video: Text-to-video generation without text-video data. *arXiv preprint arXiv:2209.14792*, 2022. 3
- [55] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning (ICML)*, 2015. 3, A1
- [56] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations (ICLR)*, 2020. 3
- [57] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. 3, A1
- [58] Sebastian Starke, He Zhang, Taku Komura, and Jun Saito. Neural state machine for character-scene interactions. *ACM Transactions on Graphics (TOG)*, 38(6):209–1, 2019. 1, 2
- [59] Andrew Szot, Alexander Clegg, Eric Undersander, Erik Wijmans, Yili Zhao, John Turner, Noah Maestre, Mustafa Mukadam, Devendra Singh Chaplot, Oleksandr Maksymets, et al. Habitat 2.0: Training home assistants to rearrange their habitat. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. 2
- [60] Omid Taheri, Nima Ghorbani, Michael J Black, and Dimitrios Tzionas. Grab: A dataset of whole-body human grasping of objects. In *European Conference on Computer Vision (ECCV)*, 2020. 2
- [61] Guy Tevet, Sigal Raab, Brian Gordon, Yonatan Shafir, Daniel Cohen-Or, and Amit H Bermano. Human motion diffusion model. *arXiv preprint arXiv:2209.14916*, 2022. 3
- [62] I Tolstikhin, O Bousquet, S Gelly, and B Schölkopf. Wasserstein auto-encoders. In *International Conference on Learning Representations (ICLR)*, 2018. 2
- [63] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017. 5

- [64] Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural Computation*, 23(7):1661–1674, 2011. 3
- [65] Ayzaan Wahid, Austin Stone, Kevin Chen, Brian Ichter, and Alexander Toshev. Learning object-conditioned exploration using distributed soft actor critic. In *Conference on Robot Learning (CoRL)*, 2021. 3
- [66] Jingbo Wang, Yu Rong, Jingyuan Liu, Sijie Yan, Dahua Lin, and Bo Dai. Towards diverse and natural scene-aware 3d human motion synthesis. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [67] Jiashun Wang, Huazhe Xu, Jingwei Xu, Sifei Liu, and Xiaolong Wang. Synthesizing long-term 3d human motion and interaction in 3d scenes. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 1, 2, 5
- [68] Jingbo Wang, Sijie Yan, Bo Dai, and Dahua Lin. Scene-aware generative network for human motion synthesis. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [69] Kai Wang, Manolis Savva, Angel X Chang, and Daniel Ritchie. Deep convolutional priors for indoor scene synthesis. *ACM Transactions on Graphics (TOG)*, 37(4):1–14, 2018. 2
- [70] Zan Wang, Yixin Chen, Tengyu Liu, Yixin Zhu, Wei Liang, and Siyuan Huang. Humanise: Language-conditioned human motion generation in 3d scenes. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. 1, 2, 5, 6
- [71] Erik Wijmans, Abhishek Kadian, Ari Morcos, Stefan Lee, Irfan Essa, Devi Parikh, Manolis Savva, and Dhruv Batra. Ddppo: Learning near-perfect pointgoal navigators from 2.5 billion frames. In *International Conference on Learning Representations (ICLR)*, 2020. 3
- [72] Jungdam Won, Deepak Gopinath, and Jessica Hodgins. Physics-based character controllers using conditional vaes. *ACM Transactions on Graphics (TOG)*, 41(4):1–12, 2022. 1, 2
- [73] Mitchell Wortsman, Kiana Ehsani, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. Learning to learn how to learn: Self-adaptive visual navigation using meta-learning. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [74] Yan Wu, Jiahao Wang, Yan Zhang, Siwei Zhang, Otmar Hilliges, Fisher Yu, and Siyu Tang. Saga: Stochastic whole-body grasping with contact. *arXiv preprint arXiv:2112.10103*, 2021. 2
- [75] Wei Yang, Xiaolong Wang, Ali Farhadi, Abhinav Gupta, and Roozbeh Mottaghi. Visual semantic navigation using scene priors. In *International Conference on Learning Representations (ICLR)*, 2019. 3
- [76] Peiyu Yu, Sirui Xie, Xiaojian Ma, Baoxiong Jia, Bo Pang, Ruiqi Gao, Yixin Zhu, Song-Chun Zhu, and Ying Nian Wu. Latent diffusion energy-based model for interpretable text modeling. In *International Conference on Machine Learning (ICML)*, 2022. 3
- [77] Ye Yuan and Kris Kitani. Dlow: Diversifying latent flows for diverse human motion prediction. In *European Conference on Computer Vision (ECCV)*, 2020. 2, 6
- [78] Mingyuan Zhang, Zhongang Cai, Liang Pan, Fangzhou Hong, Xinying Guo, Lei Yang, and Ziwei Liu. Motiandifuse: Text-driven human motion generation with diffusion model. *arXiv preprint arXiv:2208.15001*, 2022. 3
- [79] Siwei Zhang, Yan Zhang, Federica Bogo, Pollefeys Marc, and Siyu Tang. Learning motion priors for 4d human body capture in 3d scenes. In *International Conference on Computer Vision (ICCV)*, 2021. 5
- [80] Siwei Zhang, Yan Zhang, Qianli Ma, Michael J Black, and Siyu Tang. Generating person-scene interactions in 3d scenes. In *International Conference on 3D Vision (3DV)*, 2020. 2
- [81] Song-Hai Zhang, Shao-Kui Zhang, Yuan Liang, and Peter Hall. A survey of 3d indoor scene synthesis. *Journal of Computer Science and Technology*, 34(3):594–608, 2019. 2
- [82] Yan Zhang, Michael J Black, and Siyu Tang. We are more than our joints: Predicting how 3d bodies move. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 6, A2
- [83] Yan Zhang, Mohamed Hassan, Heiko Neumann, Michael J Black, and Siyu Tang. Generating 3d people in scenes without people. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 1, 2, 5, 6, A1, A2
- [84] Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip HS Torr, and Vladlen Koltun. Point transformer. In *International Conference on Computer Vision (ICCV)*, 2021. 5, A1
- [85] Kaifeng Zhao, Shaofei Wang, Yan Zhang, Thabo Beeler, and Siyu Tang. Compositional human-scene interaction synthesis with semantic control. In *European Conference on Computer Vision (ECCV)*, 2022. 2
- [86] Shengjia Zhao, Jiaming Song, and Stefano Ermon. Infovae: Information maximizing variational autoencoders. *arXiv preprint arXiv:1706.02262*, 2017. 2
- [87] Fengda Zhu, Yi Zhu, Xiaojun Chang, and Xiaodan Liang. Vision-language navigation with self-supervised auxiliary reasoning tasks. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2
- [88] Yuke Zhu, Roozbeh Mottaghi, Eric Kolve, Joseph J Lim, Abhinav Gupta, Li Fei-Fei, and Ali Farhadi. Target-driven visual navigation in indoor scenes using deep reinforcement learning. In *International Conference on Robotics and Automation (ICRA)*, 2017. 2

## A. Background for Diffusion Model

A diffusion model [19, 55] is defined by a forward process that gradually corrupts data  $\tau^0 \sim q(\tau^0)$  over  $T$  timesteps

$$q(\tau^{1:T}|\tau^0) = \prod_{t=1}^T q(\tau^t|\tau^{t-1})$$

$$q(\tau^t|\tau^{t-1}) = \mathcal{N}(\tau^t; \sqrt{1-\beta^t}\tau^{t-1}, \beta\mathbf{I})$$

and a reverse process  $p_\theta(\tau^0) = \int p_\theta(\tau^{0:T})d\tau^{1:T}$  where

$$p_\theta(\tau^{0:T}) = p(\tau^T) \prod_{t=1}^T p_\theta(\tau^{t-1}|\tau^t)$$

$$p_\theta(\tau^{t-1}|\tau^t) = \mathcal{N}(\tau^{t-1}; \mu_\theta(\tau^t, t), \Sigma_\theta(\tau^t, t)).$$

The forward process hyperparameters  $\beta^t$  are set so that  $\tau^T$  is approximately distributed according to a standard normal distribution, so  $\tau^T$  is set to a standard normal prior as well. The reverse process is trained to match the joint distribution of the forward process by optimizing the evidence lower bound (ELBO) [19, 55]. As suggested by the literature [19, 39], we can use the reverse process parametrizations as:

$$\mu_\theta(\tau^t, t) = \frac{1}{\sqrt{\alpha^t}}(\tau^t - \frac{\beta^t}{\sqrt{1-\alpha^t}}\epsilon_\theta(\tau^t, t))$$

$$\Sigma_\theta^{ii}(\tau^t, t) = \exp(\log \hat{\beta}^t + (\log \beta^t - \log \hat{\beta}^t)v_\theta^i(\tau^t, t))$$

where  $\alpha^t = 1 - \beta^t$ ,  $\hat{\alpha}^t = \sum_{s=1}^t \alpha^s$ , and  $\hat{\beta}^t = \frac{1-\hat{\alpha}^{t-1}}{1-\hat{\alpha}^t}\beta^t$ .

We can optimize modified loss instead of the ELBO to improve the sample quality, depending on whether we learn  $\Sigma$  or treat it as a fixed hyper-parameter. For the non-learned case, we use the simplified loss:

$$\mathcal{L}_{simple}(\theta) = \mathbb{E}_{t,\epsilon,\tau^0} \left[ \|\epsilon - \epsilon_\theta(\sqrt{\hat{\alpha}^t}\tau^0 + \sqrt{1-\hat{\alpha}^t}\epsilon, t)\|^2 \right]$$

$$= \mathbb{E}_{t,\epsilon,\tau^0} \left[ \|\epsilon - \epsilon_\theta(\tau^t, t)\|^2 \right]$$

It is a weighted form of the ELBO that resembles denoising score matching over multiple noise scale [19, 57].

**Conditional Diffusion Model** The goal of the conditional diffusion model is to learn a conditional distribution  $p_\theta(\tau^0|\mathbf{c})$ . We modify the diffusion model to include the condition  $\mathbf{c}$  as input to the inverse process:

$$p_\theta(\tau^{0:T}|\mathbf{c}) = p(\tau^T) \prod_{t=1}^T p_\theta(\tau^{t-1}|\tau^t, \mathbf{c})$$

$$p_\theta(\tau^{t-1}|\tau^t, \mathbf{c}) = \mathcal{N}(\tau^{t-1}; \mu_\theta(\tau^t, t, \mathbf{c}), \Sigma_\theta(\tau^t, t, \mathbf{c}))$$

## B. Model Architectures

For the tasks of human pose/motion generation in 3D scenes and path planning for 3D scene navigation, we use

the same scene encoder, *i.e.*, the PointTransformer [84] adopted from the original architecture. We pre-train the scene encoder with indoor scene semantic segmentation task on ScanNet dataset and freeze it while training SceneDiffuser. The outputs of the scene encoder are used as the key and value of the cross-attention module.

For processing the trajectory, we employ an FC layer and positional embedding to obtain the high-dimensional feature of the trajectory. We then fuse the trajectory feature with denoising timestep embedding with a ResBlock. After that, we feed the fused feature vectors to a self-attention module and use them as the query of the cross-attention module. Finally, the computed vector is fed into a feedforward layer to estimate the noise  $\epsilon$ .

For the task of dexterous grasp generation for 3D objects, we use PointNet [43] as the 3D object encoder. Before the cross-attention module, the outputs of PointNet are reshaped to  $(N_{\text{token}}, N_{\text{feat}})$ , where  $N_{\text{token}}$  refers to the number of tokens and  $N_{\text{feat}}$  refers to the dimensions of the feature.

For the task of motion planning for robot arms, we adopt PointTransformer [84] as the scene encoder, which is jointly trained from scratch with SceneDiffuser.

## C. Objective Design

For human pose and motion generation in 3D scenes, we encourage contact and non-collision between the generated human body meshes and the scene meshes. Following [83], we design optimization objective  $\varphi_o(\tau^t|\mathcal{S}) = \alpha_1\varphi_o^{\text{collision}} + \alpha_2\varphi_o^{\text{contact}}$  for pose generation and  $\varphi_o(\tau^t|\mathcal{S}) = \alpha_1\varphi_o^{\text{collision}} + \alpha_2\varphi_o^{\text{contact}} + \alpha_3\varphi_o^{\text{smoothness}}$  for motion generation.  $\alpha$  is the balancing weight.  $\varphi_o^{\text{collision}}$  minimizes the negative signed-distance values of the body mesh vertices given the negative signed distance field (SDF) of the 3D scene  $\Phi_s^-(\cdot)$ , which is formulated as

$$\varphi_o^{\text{collision}} = -\mathbb{E} \left[ |\Phi_s^-(\mathcal{M}^t)| \right], \quad (\text{A1})$$

where  $\mathcal{M}^t$  is the SMPL-X body mesh at denoising step  $t$ .  $\varphi_o^{\text{contact}}$  minimize the distance between contact body parts of the generated body mesh and the scene mesh, which is formulated as

$$\varphi_o^{\text{contact}} = - \sum_{v_c \in C(\mathcal{M}^t)} \min_{v_s \in \mathcal{S}} |v_c - v_s|, \quad (\text{A2})$$

where  $C(\cdot)$  is the operation of selecting contact part vertices from the SMPL-X body mesh according to the annotation in [15]. We design the smoothness objective to smooth the motion over time by minimizing the velocity difference of consecutive frames, which is formulated as

$$\varphi_o^{\text{smoothness}} = - \sum_{v \in \mathcal{M}^t} \sum_{i=1}^{L-2} \|v^{i+2} - 2v^{i+1} + v^i\|^2, \quad (\text{A3})$$

where  $L$  is the length of the motion sequence. We empirically set  $\alpha_1 = 1.0$ ,  $\alpha_2 = 0.02$ , and  $\alpha_3 = 0.001$ .

For dexterous grasp generation, we punish the collision between the robotic hand mesh and the object mesh. We design optimization  $\varphi_o(\tau^t|\mathcal{S}) = \varphi_o^{\text{collision}}$ .  $\varphi_o^{\text{collision}}$  is similar to Eq. (A1), where 3D scene is represented as an object and  $\mathcal{M}^t$  as the robotic hand mesh at denoising step  $t$ .

For path planning for the 3D scene navigation task, we design an optimization objective  $\varphi_o$  and  $\varphi_p$  for generating collision-free paths toward the goals. The collision-free objective maximizes the distance between the robot and the scene vertices in the robot cylinder, formulated as

$$\varphi_o = - \sum_{i=1}^L \sum_{v_s \in \mathcal{S}} \text{ReLU}(r - \text{dist}(v_s, \tau_i^t)), \quad (\text{A4})$$

where  $\text{ReLU}(x) = \max(0, x)$ ,  $r$  is the radius of the robot cylinder, and  $\text{dist}(\cdot)$  compute the Euler distance between scene vertices and robot position on the 2D plane. The planning objective  $\varphi_p$  encourages the generated paths toward the target position. In our work, we formulate it as

$$\varphi_p = \sum_{i=1}^L \exp\left(\frac{1}{\|\mathcal{G} - \tau_i^t\|_1}\right). \quad (\text{A5})$$

For robot arm motion planning, we design the planning objective  $\varphi_p$  similar to Eq. (A5). The objective is defined as

$$\varphi_p = \sum_{i=1}^L \exp\left(\frac{1}{\sum_{j=1}^N \|\mathcal{G}_j - \tau_{ij}^t\|_1}\right). \quad (\text{A6})$$

where  $N$  denotes to the number of revolute joints and  $j$  refers to  $j$ -th revolute joint.

## D. Implementation Details

### D.1. Human Pose Generation in 3D Scenes

Following prior work [83], we represent the human body with the SMPL-X model. We denote the parameters of SMPL-X in our setting as  $x_h := (t, R, \beta, \theta_b)^T \in \mathbb{R}^{79}$ , where  $t$  is the global translation in meters,  $R$  is the global orientation represented in axis-angle,  $\beta \in \mathbb{R}^{10}$  is the body shape feature, and  $\theta_b \in \mathbb{R}^{63}$  is the axis-angle representation of 21 body joints. SMPL-X can map these low-dimensional parameters into a watertight mesh with a fixed topology, enabling physical collision and contact modeling. Unlike [83] using scene depth and semantics, we directly represent the scene with a point cloud  $\mathcal{S} \in \mathbb{R}^{32768 \times 3}$ , which provides raw information about the 3D scene.

For quantitative evaluation, we randomly sample 1000 examples in each test scene to compute the diversity and physics metrics. Specifically, we separately calculate the Average Pairwise Distance (APD) and standard deviation

(std) for global translation  $t \in \mathbb{R}^3$ , the rest of local SMPL-X parameters  $(R, \beta, \theta_b)^T \in \mathbb{R}^{76}$ , and the marker-based representation [82] of generated bodies without global translation. We also report the non-collision score of the generated human bodies by calculating the proportion of the scene vertices with positive SDF to the human body and the contact score by checking whether the body contacts with the scene within a pre-defined distance threshold, *i.e.*, 0.02m.

To train SceneDiffuser, we use Adam [27] optimizer with 0.0001 as the learning rate. We use 4 NVIDIA A100 GPUs to train 100 epochs with a batch size of 128. The number of diffusion steps  $T$  in this task is set as 100. For optimization guidance sampling, we empirically set scale coefficient  $\lambda = 2.5$ .

### D.2. Human Motion Generation in 3D Scenes

For the two different settings (with and without start position) of human motion generation in 3D scenes, we represent the single-frame human body of the motion sequence as the same as the pose generation. To collect training data, we clip the motion sequences in the PROX dataset into motion segments with a fixed duration, *i.e.*, 60 frames. We use the same evaluation metrics as pose generation and report the average values over motion sequence as the motion generation performance measure. In this task, the optimizer is Adam, and the learning rate is 0.0001. We use 4 NVIDIA A100 GPUs to train 300 epochs with 200 diffusion steps and 128 batch size. For optimization guidance sampling, we empirically set scale coefficient  $\lambda = 2.5$ .

### D.3. Dexterous Grasp Generation for 3D Objects

We use Shadowhand as our dexterous robotic hand and denote qpos as  $q := (t, R, \theta) \in \mathbb{R}^{33}$ , where  $t \in \mathbb{R}^3$  and  $R \in \mathbb{R}^6$  represent the global translation and orientation respectively,  $\theta \in \mathbb{R}^{24}$  describes the rotation angles of the revolute joints. We split the MultiDex [31] into 48 seen objects and 10 unseen objects for training and testing.

For each grasp, we apply  $0.5\text{ms}^{-2}$  acceleration to the object along  $\pm xyz$  directions, and the grasping is successful if the movements of the object are all within 2cm. For the diversity, we first capture the mean  $\mu_i$  and the standard deviation  $\sigma_i$  of  $i$  revolute joint in the training data grasping pose. We define the mean pose as  $\mu_q := (\mu_1, \mu_2, \dots, \mu_{24}) \in \mathbb{R}^{24}$  and the standard deviation pose as  $\sigma_q := (\sigma_1, \sigma_2, \dots, \sigma_{24}) \in \mathbb{R}^{24}$ . We report the success rate of generated poses that lie at the  $k$  standard deviation level, which means these poses  $q$  satisfy the constraint as  $\mu_q - k\sigma \leq q \leq \mu_q + k\sigma$ . For the depth collision computation, we sample the surface points  $\mathcal{H} \in \mathbb{R}^{3200 \times 3}$  on the ShadowHand related to the pose  $q$  and the surface points with normal  $\mathcal{O} \in \mathbb{R}^{4096 \times 6}$  on the object. We compute the collision for ShadowHand surface to the object and report the depth collision among  $\mathcal{H}$  to show the quality of generated poses.

To train SceneDiffuser on this task, we use Adam optimizer, set the learning rate as 0.0001, and use 1 NVIDIA A100 GPU to train 2100 epochs with 64 batch size. For optimization guidance sampling, we empirically set scale coefficient  $\lambda = 1.0$ .

#### D.4. Path Planning for 3D Scene Navigation

In this task, we consider 3D navigation in realistic scenes, where the goal is to plan plausible trajectories for a physical robot from the given start position  $\hat{s}_0$  to the given target position  $\mathcal{G}$  in a furnished 3D indoor scene  $\mathcal{S}$ . We represent the hallucinated physical robot as a cylinder to simulate physically plausible trajectories which are collision-free in the 3D scene. The robot can move in all directions within a distance in each step without height change. We set the maximum moving distance as 0.08m, the robot radius as 0.08m, and the robot height as infinite, which means the robot can only move on the floor that is not occupied.

To construct room-level realistic scenarios for path planning, we manually select 61 indoor scenes from ScanNet [9], as shown in Fig. A1. We annotate these scenes with spatially dense and physically plausible navigation graphs and collect about 6.3k trajectories by searching the shortest paths between the randomly selected start and target graph nodes. As the distance between nodes may be too long for a robot to move in one step, we refined the trajectories according to the maximum moving distance. These trajectories have an average step of 60.0, a minimal step of 32, and a maximum step of 120. We use 4.7k trajectories in 46 scenes as the training data and the rest 1.6k trajectories in 15 unseen scenes for evaluation. We set the maximum number of planning steps as 150.

During training, we set the fixed trajectory horizon as 32. We use 4 NVIDIA A100 GPUs to train 50 epochs with 100 diffusion steps and a batch size of 128. The optimizer is Adam, and the learning rate is 0.0001. During inference, we empirically set the scale coefficient of optimization guidance as 1.0 and the scale coefficient of planning guidance as 0.2.

#### D.5. Motion Planning for Robot Arm

We use the Franka Emika with seven revolute joints as the robot arm and randomly generate cluttered tabletop scenes with primitives following specific heuristics. For each scene, we position the robot arm at the center of the table and use moveit2 motion planner [48] to synthesize trajectories constrained by a pair of start and goal poses of the end effector. We collected 19,800 collision-free trajectories over 200 clustered scenes.

During inference, we execute the planned motions of SceneDiffuser in IsaacGym [35]. We consider the planning is successful if our robot arm reaches the goal pose by a certain  $L_2$  norm distance (e.g., 0.2) in the space of revolute joints. Note that the simulation can not run infinitely;

therefore, we set a limited number of simulation steps (e.g., 300). For the efficiency evaluation, we capture the average number of simulation steps.

To train SceneDiffuser on this task, we use Adam optimizer, set the learning rate as 0.0001, and use 4 NVIDIA A100 GPUs to train 200 epochs with 128 batch size. We empirically set the scale coefficient of planning guidance as 0.2 during inference.

#### D.6. Scaling Factor for the Guidance

Similar to Ho *et al.* [19], we notice that the parameter  $\Sigma$  in Eq. (9) decreases as the denoising step  $t$  decreases, which gradually weakens the guidance during the denoising process. Instead of using a constant as the scaling factor, we empirically schedule the scaling factor by dividing it by  $\Sigma$ . It reformulates Eq. (9) as

$$p(\tau^{t-1} | \tau^t, \mathcal{S}, \mathcal{G}) \approx \mathcal{N}(\tau^{t-1}; \mu + \lambda \mathbf{g}, \Sigma)$$

### E. Additional Ablative Experiments

We ablate different model architectures, including the scene encoder and noise prediction module in SceneDiffuser, diffusion steps and scale coefficient in the optimizer of dexterous grasp generation task, and fixed frames and planning objectives of path planning for 3D scene navigation task.

#### E.1. Model Architecture

As shown in Tab. A1, we study how different scene model influences the dexterous grasp generation results. We use PointNet [43] and PointNet++ [44] as different scene models to extract the object feature. For more diversity evaluation, we capture the mean standard deviation among all revolute joints of the robotic hand qpos. We find that the global feature extracting from PointNet makes it easier for the model to learn a mean pose to obtain a higher grasping success rate. In contrast, the local feature extracting from PointNet++ makes the generated grasp pose more diverse.

Table A1. Ablation on different scene encoder.

Scene Encoder	Succ. Rate (%) $\uparrow$			Div. (rad.) $\uparrow$	Coll. (mm) $\downarrow$
	$\sigma$	$2\sigma$	all		
PointNet (w/o opt.)	70.65	71.25	71.25	0.0718	17.34
PointNet (w/ opt.)	71.27	70.32	69.84	0.0838	14.61
PointNet++ (w/o opt.)	56.47	66.29	66.25	0.1568	18.53
PointNet++ (w/ opt.)	64.33	60.51	59.53	0.1670	14.37

As shown in Tab. A2, we ablate the module for noise prediction. We compare the design of cross-attention and self-attention for processing the condition and input. Cross-attention indicates learning query from the input  $\tau_t$  and learning key and value from the scene condition  $\mathcal{S}$ . Self-attention indicates concatenating  $\tau_t$  and scene features  $\mathcal{S}$

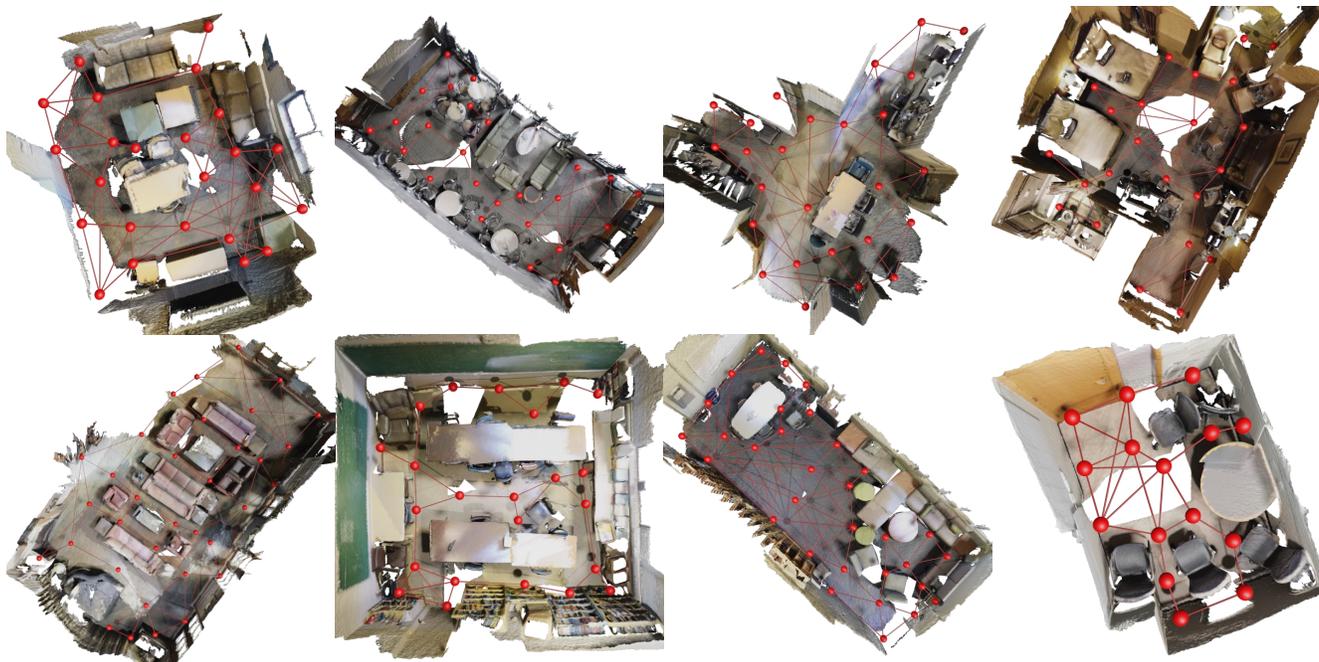


Figure A1. Scenes and corresponding navigation graphs for path planning. The selected scenes have various regions, diverse room types, and sufficient layout complexity.

and learning with self-attention. Through our experiments, we find that with self-attention, the model learns better to capture the joint distribution of input and condition. This leads to a slightly lower diversity but better generation quality and success rate.

Table A2. Ablation on different model architecture.

Epsilon Model	Succ. Rate (%) $\uparrow$			Div. (rad.) $\uparrow$	Coll. (mm) $\downarrow$
	$\sigma$	$2\sigma$	all		
CrossAttn. (w/o opt.)	70.65	71.25	71.25	0.0718	17.34
CrossAttn. (w/ opt.)	71.27	70.32	69.84	0.0838	14.61
SelfAttn. (w/o opt.)	74.27	75.94	75.94	0.0535	16.49
SelfAttn. (w/ opt.)	72.01	71.56	71.09	0.0605	13.94

## E.2. Diffusion Steps

We study different diffusion steps  $T$  in Tab. A3, where we use PointNet++ as the scene encoder with cross-attention design. We report the success rate, diversity, and depth collision of sampling results in the test set under different diffusion steps, ranging from 30 to 1000.  $T$  balance the diversity and success rate in dexterous grasp generation, where  $T = 30$  leads to the best diversity of generated grasp pose and  $T = 1000$  leads to the best *all* success rate.

## E.3. Scale Coefficient

Among different time steps  $T$ , we ablate scale coefficient  $\lambda$  of the optimization guidance in dexterous grasp genera-

tion in Tab. A3, ranging from 0.0 (denoted as w/o in the table) to 1.0. Through extensive experiments, we observe that, in general, the  $\alpha$  trade off the depth collision and grasp success rate. A larger  $\alpha$  value leads to fewer collisions and draws the grasp pose away from the object simultaneously, which loses the grasp stability and lowers the success rate.

We also ablate the scale coefficient of the planner in path planning for 3D scene navigation, as shown in Tab. A4. Too small or too large scale coefficients both lead to a performance drop. It is due to that a small value cannot provide sufficient guidance. In contrast, a large value diminishes trajectory diversity with strong guidance, preventing it from escaping obstacles and dead-ends.

## E.4. Fixed Frames for Planning

Since we formulate the planning algorithm as inpainting, we also ablate the number of the fixed frame in it. In path planning for 3D scene navigation, we train the SceneDiffuser with a trajectory length of 32. Therefore, we compare the settings of fixing the first 1, 7, 15, 23, and 31 frames for inpainting during the denoising process. The results in Tab. A4 show that the model achieves the best performance while fixing the first 15 frames.

## E.5. Planning Objectives

To explore the influence of different planning objectives, we design the following four planning objectives and compare them with Eq. (A5).

Table A3. Ablation on diffusion steps and scale coefficient.

Time Steps	Optimizer Scale	Succ. Rate (%) $\uparrow$				Diversity (rad.) $\uparrow$	Depth Collision (mm) $\downarrow$
		$\sigma$	$2\sigma$	$3\sigma$	all		
30	w/o	0.00	60.01	50.94	48.13	0.3418	21.19
30	0.1	0.00	58.72	54.90	51.09	0.3415	19.96
30	0.5	0.00	64.24	51.63	47.81	0.3397	17.41
30	1.0	0.00	60.41	48.76	43.59	0.3393	16.05
100	w/o	0.00	66.62	60.12	58.91	0.2865	19.07
100	0.1	0.00	66.54	60.60	59.53	0.2836	17.55
100	0.5	0.00	61.23	56.71	53.75	0.2898	14.63
100	1.0	0.00	56.79	53.13	48.91	0.2920	14.53
500	w/o	75.00	67.50	67.34	67.34	0.1753	19.29
500	0.1	68.56	65.19	65.00	65.00	0.1733	17.68
500	0.5	62.83	60.25	58.94	58.75	0.1814	15.12
500	1.0	62.21	57.76	55.17	54.37	0.1872	14.36
1000	w/o	56.47	66.29	66.26	66.25	0.1568	18.53
1000	0.1	73.24	71.43	71.04	71.09	0.1572	16.88
1000	0.5	70.18	65.99	65.55	65.62	0.1611	14.37
1000	1.0	64.33	60.51	59.61	59.53	0.1670	14.37

Table A4. Ablation on different inpainting horizons and scale coefficients of the planning guidance.

Fixed Frames	Planner Scale	Succ. Rate (%) $\uparrow$	Planning Steps $\downarrow$
1	0.2	31.25	135.14
7	0.2	65.50	104.30
15	0.2	73.75	90.38
23	0.2	73.25	87.49
31	0.1	53.50	106.23
	0.2	62.37	97.02
	0.3	56.81	101.54
	0.4	50.94	105.11

- We compute the L1 distance between the last frame of the denoised trajectory and the target position, *i.e.*,

$$\varphi_p = -\|\mathcal{G} - \tau_L^t\|_1. \quad (\text{A7})$$

- We summarize the L1 distance between all frames of the denoised trajectory and the target position, *i.e.*,

$$\varphi_p = -\sum_{i=1}^L \|\mathcal{G} - \tau_i^t\|_1. \quad (\text{A8})$$

- Similar to Eq. (A5), we only consider the last frame of the denoised trajectory, *i.e.*,

$$\varphi_p = \exp\left(\frac{1}{\|\mathcal{G} - \tau_L^t\|_1}\right). \quad (\text{A9})$$

- We compute the L1 distance between the target position and the frame closest to the target, *i.e.*,

$$\varphi_p = -\min_i \|\mathcal{G} - \tau_i^t\|_1. \quad (\text{A10})$$

The planning results in Tab. A5 indicate that encouraging all frames of the denoised trajectory to reach the target position surpasses considering only one frame. Besides, directly using L1 distance tends to achieve a better performance than additionally applying the exponential function.

Table A5. Ablation on different planning objectives.

Objective	Succ. Rate (%) $\uparrow$	Planning Steps $\downarrow$
$\varphi_p = -\ \mathcal{G} - \tau_L^t\ _1$	57.06	116.22
$\varphi_p = -\sum_{i=1}^L \ \mathcal{G} - \tau_i^t\ _1$	75.69	88.02
$\varphi_p = \exp\left(\frac{1}{\ \mathcal{G} - \tau_L^t\ _1}\right)$	34.31	131.74
$\varphi_p = \sum_{i=1}^L \exp\left(\frac{1}{\ \mathcal{G} - \tau_i^t\ _1}\right)$	73.75	90.38
$\varphi_p = -\min_i \ \mathcal{G} - \tau_i^t\ _1$	56.00	109.02

## F. Trainable Optimization and Planning

As shown in Alg. 1, we can optionally train the optimization and planning process with observed trajectories. To verify its efficacy, we optimize the trainable scaling factor  $\lambda$  of the optimization guidance in pose generation and path planning tasks. Specifically, we use a small MLP model to map the timestep embedding of each step into a scalar, *i.e.*, the scaling factor. During training, we only optimize the MLP while fixing the pre-trained diffusion model. We plot the learned scaling factor varying with the denoising step from 100 to 1, as shown in Fig. A2. We observe that the scaling factor of the denoising process at the beginning is much smaller than at the end. We speculate that the target signal at the beginning of the denoising process is mostly noise so a large scaling factor cannot optimize it properly. The scaling factor decrease in the last several steps may be because this can alleviate excessive guidance and balance the guidance from other modules, such as the planner.

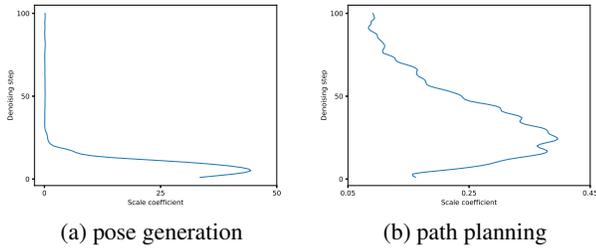


Figure A2. Trainable scaling factor varying with the denoising step.

## G. More Qualitative Results

**Pose Generation in 3D Scenes** We show more qualitative results in Fig. A3.

**Motion Generation in 3D Scenes** We provide more sampled human motions from the same start pose in other scenes, as shown in Fig. A4. Please refer to the supplemental demo video for better visualization with rendered animations.

**Path Planning for 3D Scene Navigation** Fig. A5 shows some qualitative results of path planning for 3D scene navigation.

**Dexterous Grasp Generation for 3D Objects** We show more qualitative results in Fig. A6. Note that the objects are unseen during training time.

**Motion Planning for Robot Arm** We render the planning results into animations for visualization. Please refer to the supplemental demo video for the qualitative results.

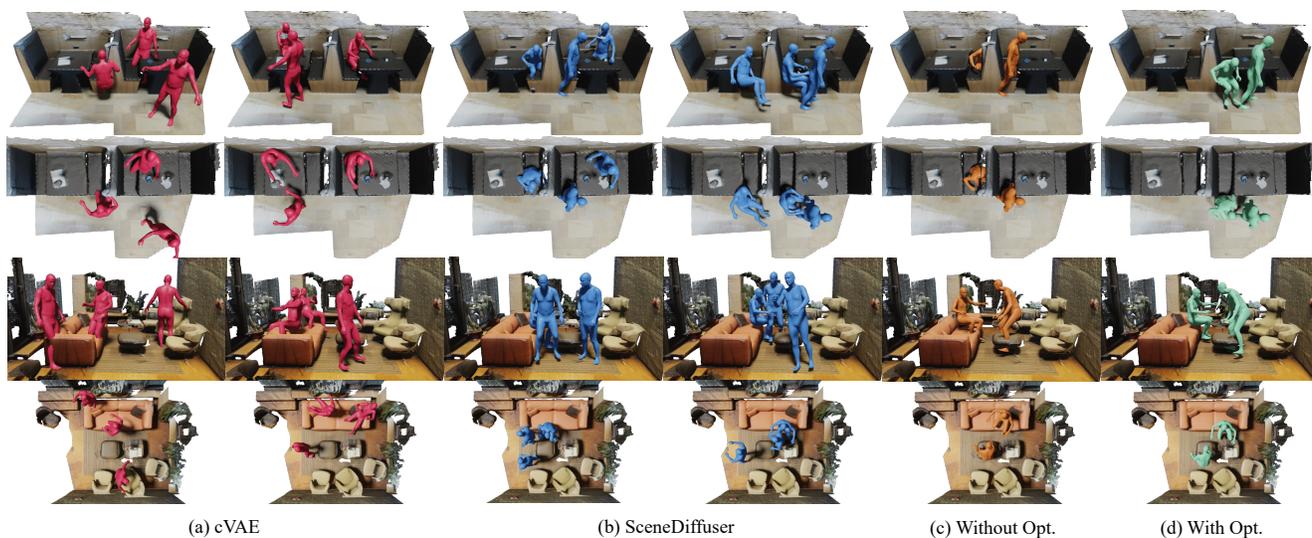


Figure A3. More qualitative results of pose generation in 3D scenes.

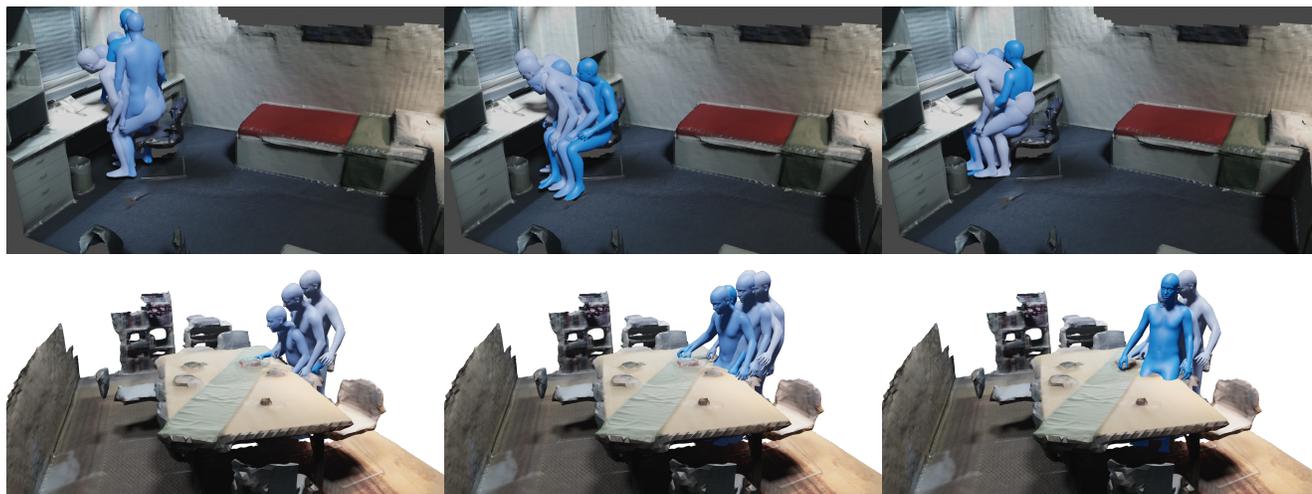


Figure A4. More qualitative results of motion generation in 3D scenes.

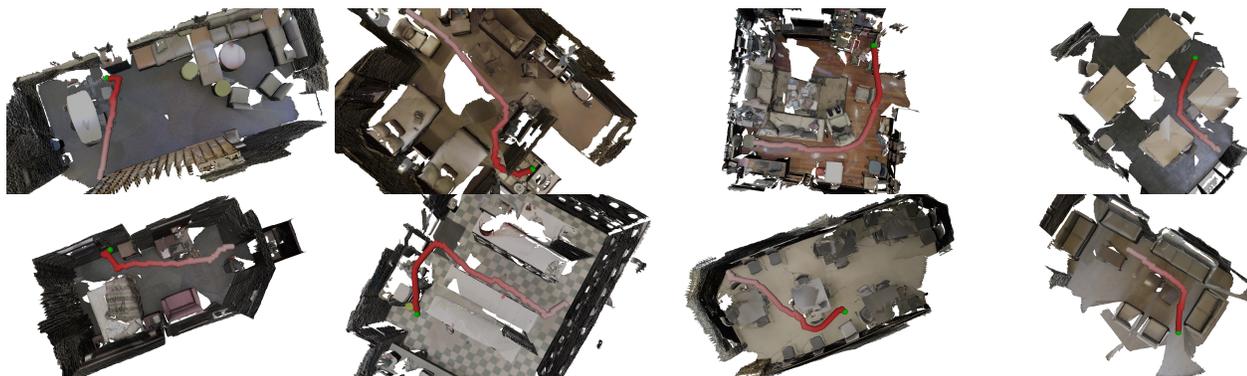


Figure A5. Qualitative results of path planning for 3D scene navigation. The red balls represent the planning result, starting with the lightest red ball and ending with the darkest red ball. The green ball indicates the target position.

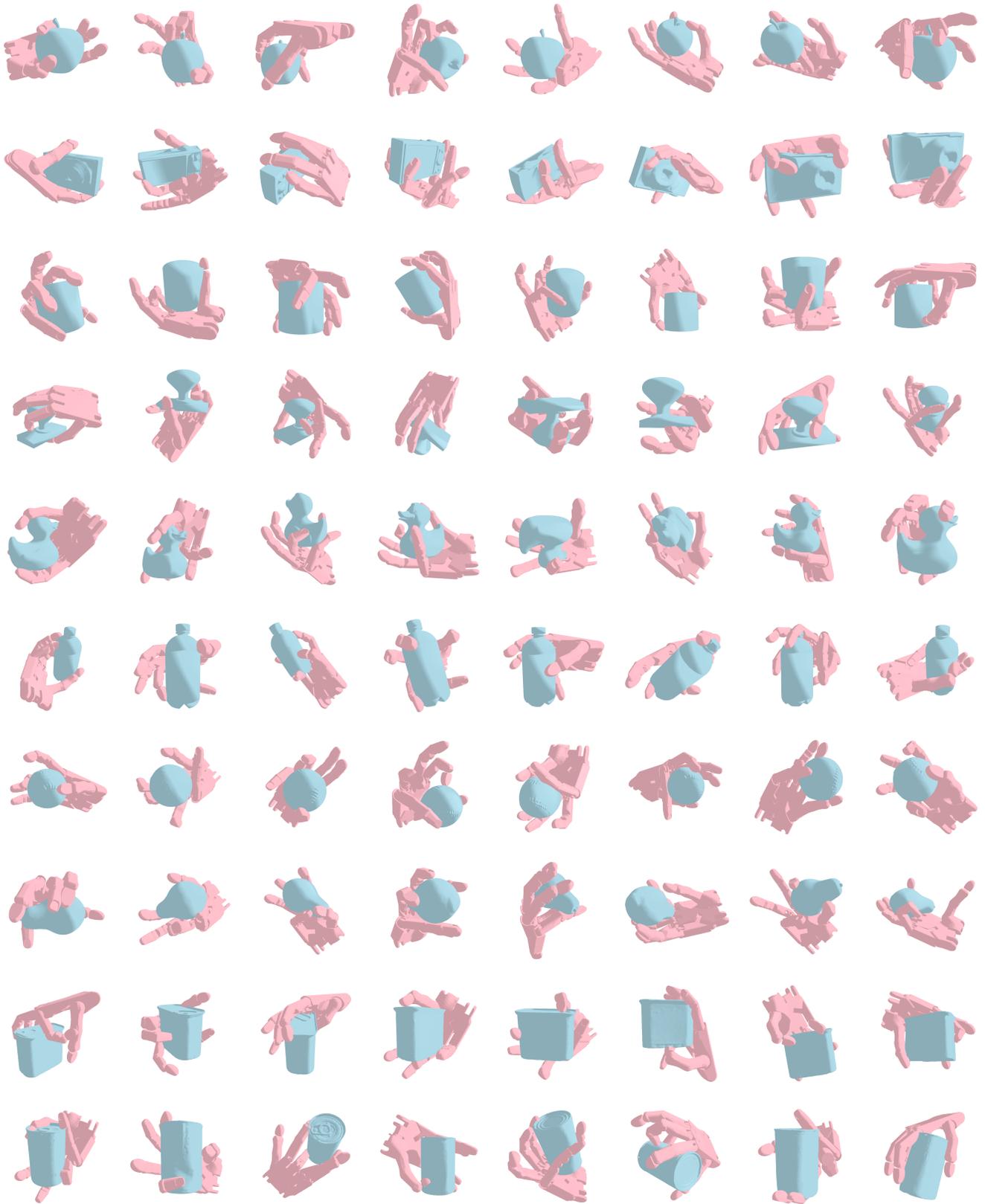


Figure A6. More qualitative results of dexterous grasp pose generation for 3D object.