



## SceneDiffuser



- human pose generation
- human motion generation
- dexterous grasp generation
- path planning for 3D navigation with goals
- motion planning for robot arms

*SceneDiffuser* is a conditional generative model for 3D scene understanding.

It is applicable to various scene-conditioned 3D tasks.

## Long-standing Goals for 3D Scene Understanding



Scene-aware Generation Physics-based Optimization Goal-oriented Planning

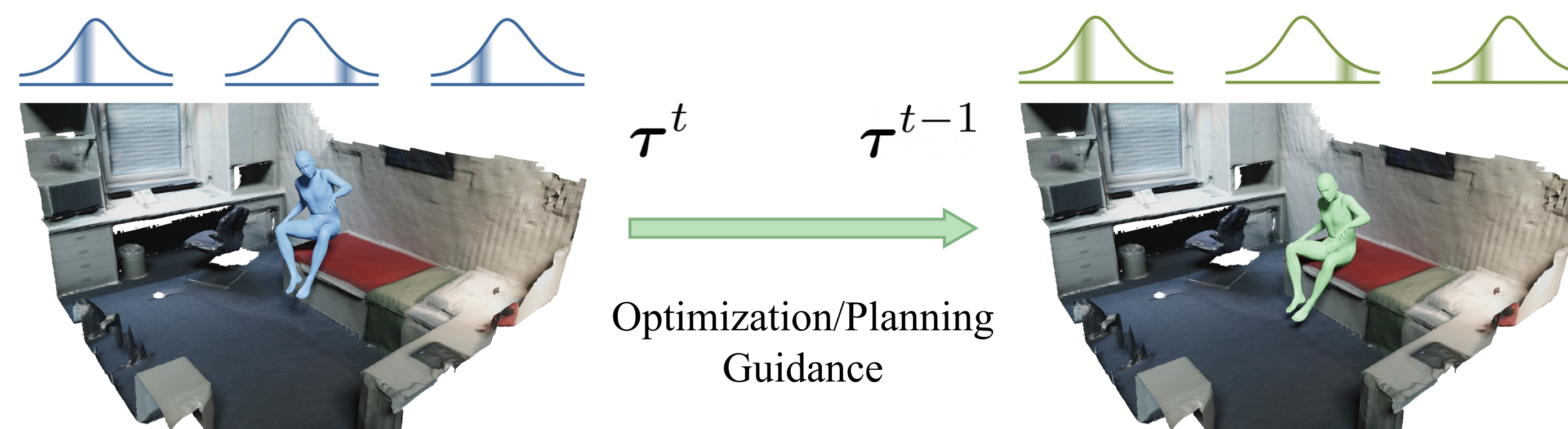
## Two Fundamental Limitations

- Lack of *powerful* generative model
- Lack of *unified* framework

## Contribution

- ✓ We propose the **SceneDiffuser** as a general conditional generative model for *generation*, *optimization*, and *planning* in 3D scenes.
- ✓ **SceneDiffuser** is intrinsically *scene-aware*, *physics-based*, and *goal-oriented*, applicable to various scene-conditioned 3D tasks.
- ✓ We demonstrate that the **SceneDiffuser** outperforms previous models by a *large margin* on **five** scene understanding tasks, establishing its **efficacy** and **flexibility**.

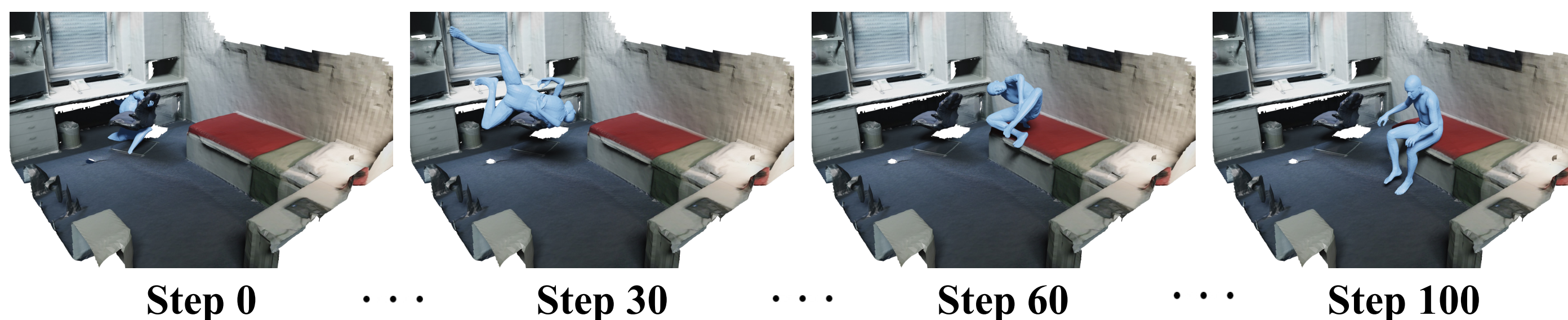
## Sampling with Optimization/Planning Guidance



$$\mu = \mu_{\theta}(\tau^t, t, \mathcal{S}), \Sigma = \Sigma_{\theta}(\tau^t, t, \mathcal{S})$$

$$\tau^{t-1} = \mathcal{N}(\tau^{t-1}; \mu + \lambda \Sigma \nabla_{\tau^t} (\mathcal{J}(\tau^t | \mathcal{S}, \mathcal{G})) |_{\tau^t = \mu}, \Sigma)$$

## Denoising Process with Guidance



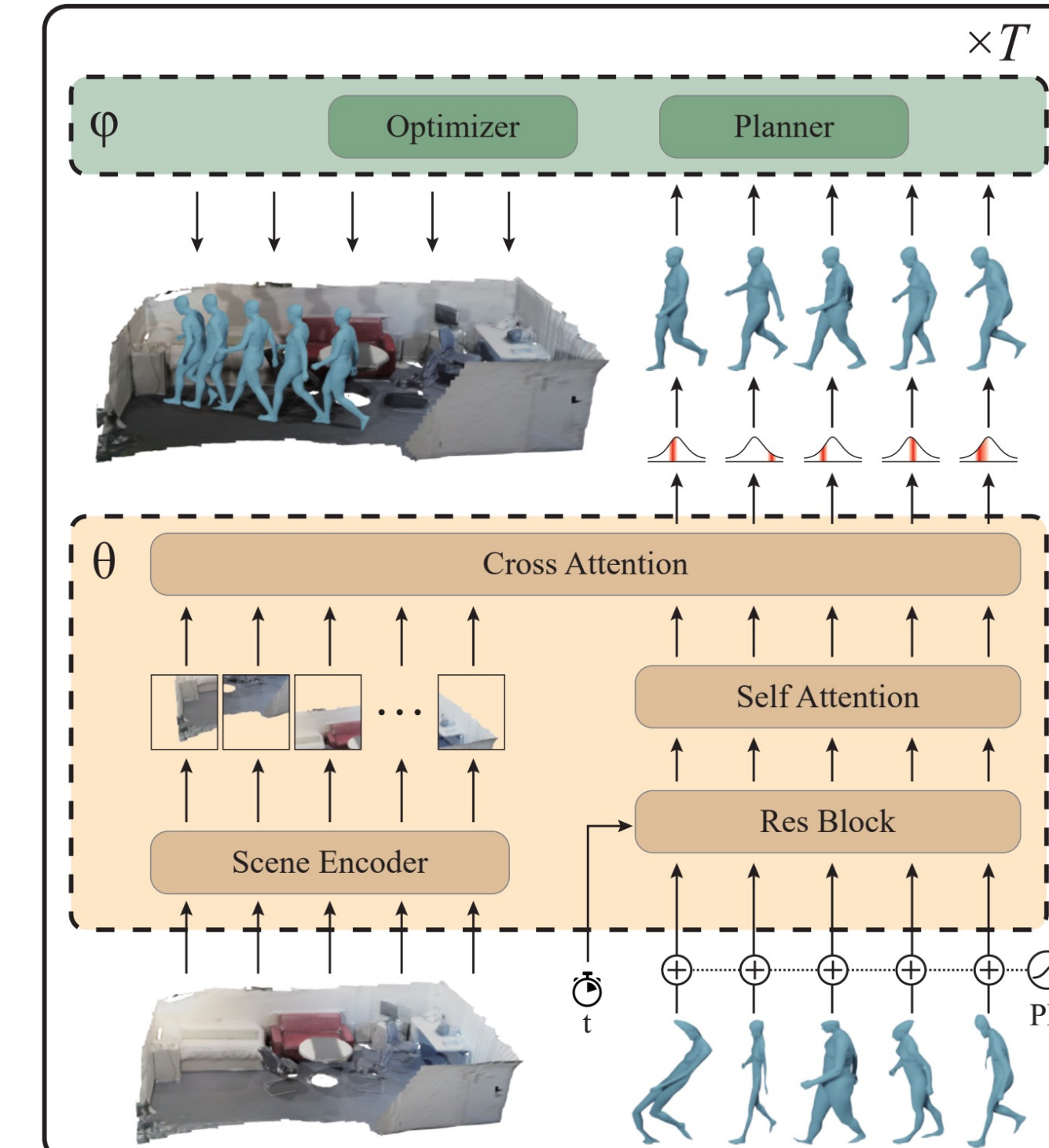
## Sampling Algorithm and Model Architecture

**Algorithm 2:** Sampling SceneDiffuser for generation, optimization, and planning

```

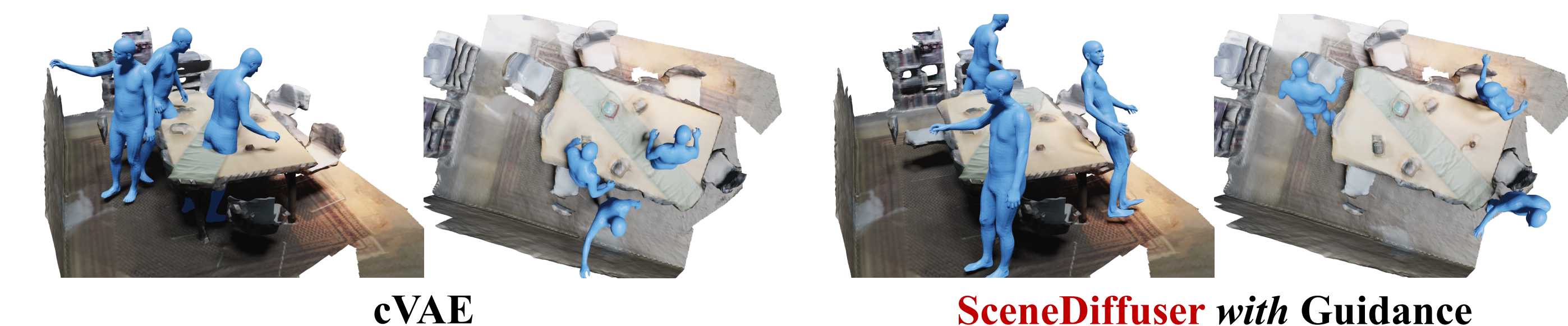
Modules: Model  $p_{\theta}(\cdot | \mathcal{S})$ , optimization objective  $\varphi_o(\cdot | \mathcal{S})$ , and planner objective  $\varphi_p(\cdot | \mathcal{S}, \mathcal{G})$ 
1 // one-step guided sampling
2 function sample( $\tau^t, \mathcal{J}$ ):
3    $\mu = \mu_{\theta}(\tau^t, t, \mathcal{S}), \Sigma = \Sigma_{\theta}(\tau^t, t, \mathcal{S})$ 
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  $\tau^T \sim \mathcal{N}(0, \mathbf{I})$ 
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  $\hat{s}_0$ , initial plan  $\tau_0^T \sim \mathcal{N}(0, \mathbf{I})$ 
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      $\tau_i^{t-1}[0:i] = \hat{s}_{0:i}$ 
18   Act  $\hat{a}_{i-1}$  to reach  $\hat{s}_i = \tau_0^0[i]$ ,  $\hat{s}_{0:i} = \hat{s}_{0:i-1} \cup \hat{s}_i$ 
19   Increment planning step  $i = i + 1$ 

```

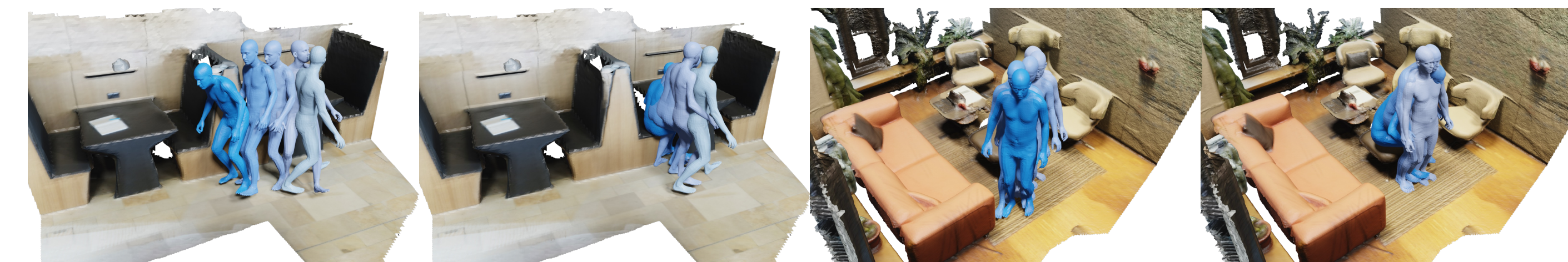


## Tasks and Results

### Task 1: Human Pose Generation

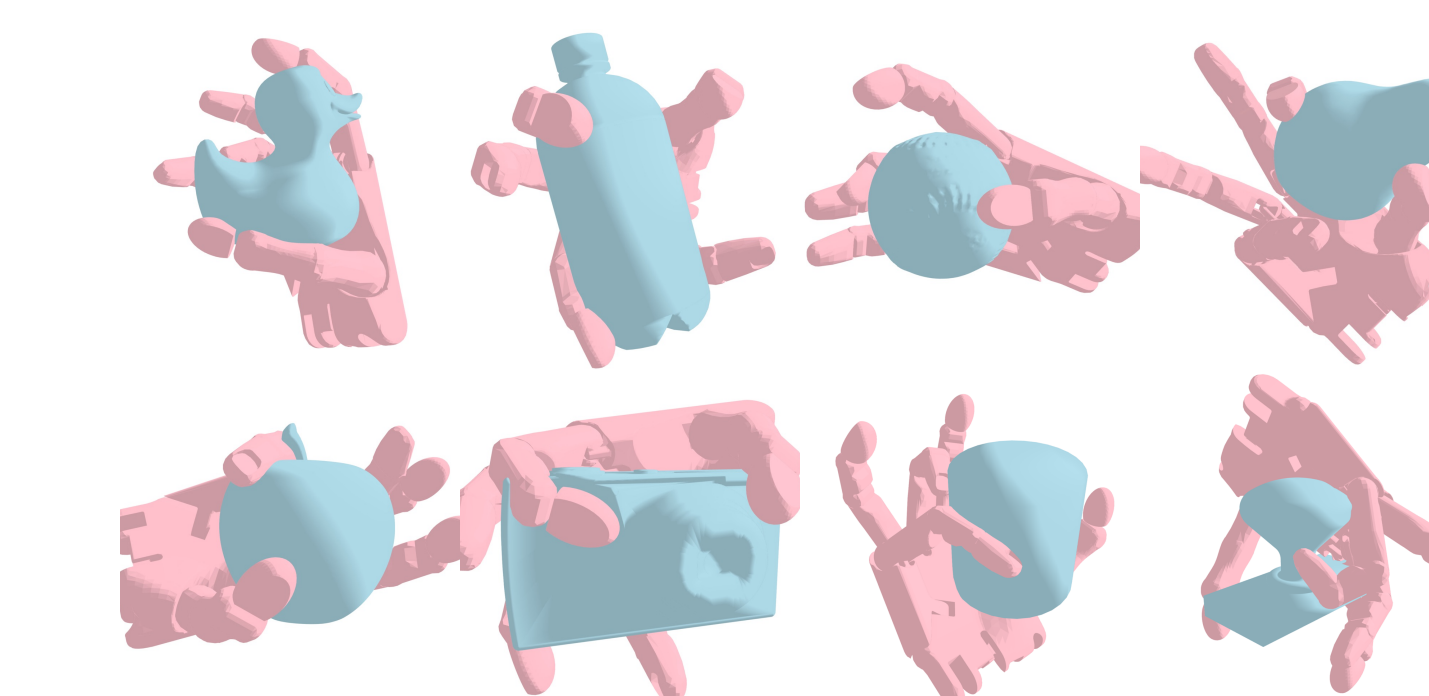


### Task 2: Human Motion Generation



**SceneDiffuser** generates diverse motions (e.g., “sit,” “walk”) from the same start position in unseen 3D scenes.

### Task 3: Dexterous Grasp Generation



**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 (%) <sup>↑</sup>			depth coll. (mm) <sup>↓</sup>
	$\sigma$	$2\sigma$	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>

### Task 4: Path Planning for Navigation

### Task 5: Motion Planning for Robot Arms

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

task	model	succ. rate (%) <sup>↑</sup>	planning steps <sup>↓</sup>
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

Please refer to our paper for more quantitative and qualitative results.