

# Learning 4D Embodied World Models

Anonymous CVPR submission

Paper ID 11899

## Abstract

In this paper, we present a 4D embodied world model, which takes in an image observation and language instruction as input and predicts a 4D dynamic scene showing how the scene will change as the embodied agent performs actions. In contrast to previously learned world models which typically generate 2D videos, our 4D model provides detailed 3D information on precise configurations and shapes of objects in a scene over time. This allows us to effectively learn accurate inverse dynamic models for an embodied agent to execute a policy for interacting with the environment. To construct a dataset to train such 4D world models, we first annotate large-scale existing video robotics datasets using pre-trained depth and normal prediction models to construct 4D-consistent scenes of each video. To efficiently learn generative models on this 4D data, we propose to train a video generative model on this annotated dataset, which jointly predicts RGB-DN (RGB, Depth, and Normal) for each video. We then present an algorithm to directly convert generated RGB, Depth, and Normal images into a high-quality dynamic 4D mesh of the world. We illustrate how this enables us to predict high-quality meshes consistent across both time and space from embodied scenarios, render novel views for embodied scenes, and construct policies that substantially outperform those from prior 2D and 3D models of the world. Our code, model, and dataset will be made publicly available.

## 1. Introduction

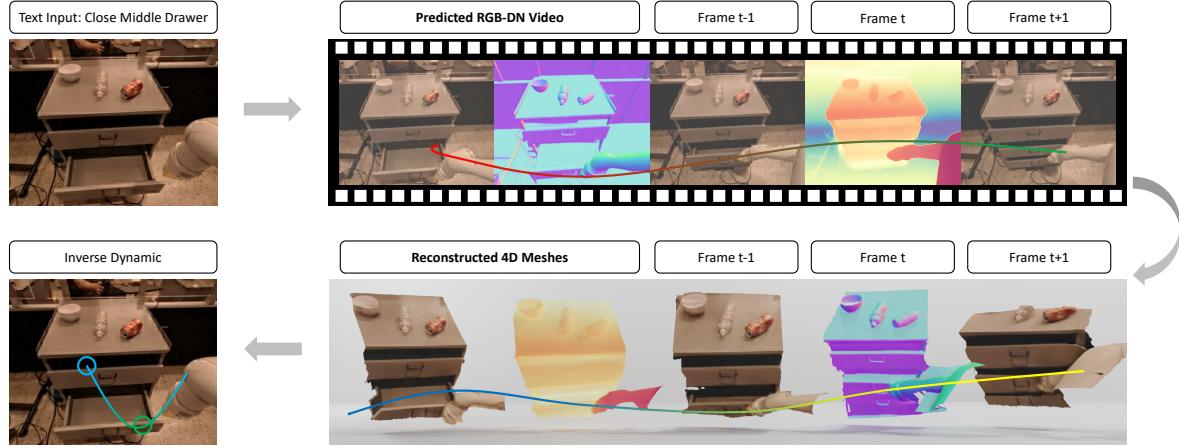
Learned world models [20, 61, 64, 70], which simulate environmental dynamics, play a crucial role in enabling embodied intelligent agents. Such models enable flexible policy synthesis [14, 41], data simulation and generation [64, 73], and long-horizon planning [13, 28, 68]. However, while the physical world is three-dimensional in nature, existing world models operate in the space of 2D pixels. This limitation leads to an incomplete representation of spatial relationships, impeding tasks that require precise depth and pose information. For instance, without accurate depth and 6-DoF pose

estimations, robotic systems struggle to determine the exact position and orientation of objects. Furthermore, existing 2D models can produce unrealistic results, such as inconsistent object sizes and shapes across time steps, which limits their use in data-driven simulations and robust policy learning.

In this paper, we explore how we can instead learn a 4D embodied world model, which directly simulates the dynamics of a 3D world. This approach allows us to generate realistic 3D interactions, such as grasping objects or opening drawers, with a level of detail that traditional 2D-based models cannot achieve. By modeling spatial and temporal dimensions, our model provides the depth and pose information essential for robotic manipulation.

However, the task of learning a 4D embodied world model is challenging as the dynamics of the world are extremely computationally expensive to train and learn, requiring models to generate outputs in three-dimensional space and time. To efficiently represent and predict the dynamics of the world, we propose a substantially more lightweight representation of the 4D world, consisting of predicting a sequence of RGB, depth, and normal maps of the scene. This combined representation accurately captures the appearance, geometry, and surface of a scene while being substantially lower dimensional than explicitly predicting world dynamics. Furthermore, such a representation shares substantial similarities to existing video space models, allowing us to directly use the generative capabilities and architecture improvements of existing video space models to effectively construct our 4D world model.

Given this intermediate representation, we present an efficient algorithm to reconstruct accurate 4D scenes from generated maps. For each frame, we use a combination of both depth and normal prediction to integrate a smooth 3D surface of the scene. We then use optical flow between generated frames to distinguish between background and dynamic regions in the reconstructed 3D scene across frames and add a loss function to reconstructions to enforce consistency across scenes over time. We find this enables us to construct fidelity 4D generated meshes for the scene suitable for downstream tasks such as policy prediction for the robot (Figure 1).



**Figure 1. 4D Embodied World Models.** Our approach gets an input image and instructions and predicts RGB-D-N (RGB, Depth, Normal) maps. We propose a normal integration method that constructs a high-quality 4D mesh of the interaction from these predictions.

079 A key challenge for training a 4D world model is a lack  
 080 of access to existing large-scale datasets with existing 4D  
 081 annotations, or the high-quality image, depth, and normal  
 082 annotations needed to train our approach. To construct such  
 083 data, in principle, we can use pre-trained depth and normal  
 084 estimators [30, 63] to obtain such estimates from video data,  
 085 but these estimators are typically limited to predicting relative  
 086 value maps from individual frames. As a result, when  
 087 a scene changes, these estimators will produce depth and  
 088 normal maps with inconsistencies across time. To tackle this  
 089 challenge, we develop a data collection pipeline that lever-  
 090 ages optical flow between frames in a video to enforce con-  
 091 sistency between generated depth and normal maps across  
 092 timesteps. In particular, we use optical flow to guide a depth  
 093 and normal diffusion model across frames in a video, which  
 094 we find is sufficient to ensure consistent depth and normal  
 095 predictions across all timesteps without the need for expen-  
 096 sive ground-truth annotations, facilitating the training of our  
 097 world model on a large scale.

098 Overall, our paper has the following contributions: **(1)** We introduce a 4D embodied world model, and present an  
 099 efficient representation of this model, in the form of RGB,  
 100 depth, and normal maps, and illustrate how this represen-  
 101 tation can be used to construct a full 4D scene. **(2)** We  
 102 present a pipeline to automatically extract 4D world model  
 103 data from existing robot video datasets, leveraging an opti-  
 104 cal flow-guided depth and normal diffusion model. **(3)** We  
 105 illustrate the approach’s efficacy in generating consistent 4D  
 106 meshes across different environments, substantially outper-  
 107 forming other baselines, and showing its downstream use as  
 108 an effective policy.

## 110 2. Related Work

111 **Embodied Foundation Models** A flurry of recent work  
 112 has focused on constructing foundation models for general

113 purpose agents [17, 65]. One line of work has focused on  
 114 constructing multimodal language models that operate over  
 115 images [12, 19, 29, 40, 50, 60, 67] as well as 3D inputs [25,  
 116 26] and output text describing the actions of an agent. Other  
 117 works have focused on the construction of vision-language-  
 118 action (VLA) models that directly output action tokens [6,  
 119 32, 69]. Both of the previous approaches aim to construct  
 120 foundation model policies (over text or continuous actions).  
 121 In contrast, our work aims to instead construct a foundation  
 122 4D world model for embodied agents, which can then be  
 123 used for downstream applications such as planning [13, 68]  
 124 or policy synthesis [14, 41].

125 **Learning World Models** Learning dynamics model of  
 126 the world given control inputs has been a long-standing  
 127 challenge in system identification [42], model-based rein-  
 128 forcement learning [54], and optimal control [2, 74]. A large  
 129 body of work focused on learning world models in the low  
 130 dimensional state space [1, 16, 38], which while being effi-  
 131 cient to learn, is difficult to generalize across many environ-  
 132 ments. Other works have explored how world models may  
 133 be learned over pixel-space images [10, 11, 20, 21, 44], but  
 134 such models are trained on simple game environments. With  
 135 advances in generative modeling, a large flurry of recent  
 136 research has focused on using video models as foundation  
 137 world models [7, 47, 61, 64, 70, 72] but such models operate  
 138 over the space of 2D pixels which does not fully simulate  
 139 the 3D world. Most similar to our work, in Zhen et al. [69],  
 140 a world model over 3D inputs is learned. In contrast to this  
 141 work, our world model directly captures the dynamics of 3D  
 142 scenes using a compact representation of RGB-DN video.

## 143 3. Learning 4D Embodied World Models

144 We introduce the 4D Embodied World Model, which predicts  
 145 future RGB, depth, and normal maps based on a given input  
 146 image and text. We leverage a pre-trained video diffusion

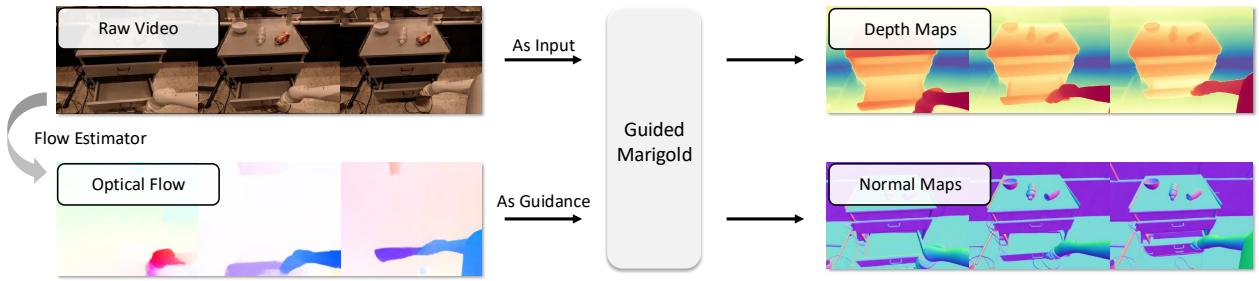


Figure 2. **Data Annotation Visualization.** An illustration for the data annotation pipeline.

model as the backbone to predict this rich set of 2D geometric information. Then, we propose an efficient method to convert the RGB-DN video into 4D scenes.

### 3.1. 4D Embodied Video Data Annotation

Learning 4D embodied world models requires large-scale 4D datasets, which are expensive to collect in the real world. In this section, we present a data annotation pipeline that enables us to automatically construct 4D datasets from existing video datasets. For an illustration of this process, see Fig. 2.

Given an input video  $\mathcal{V}$ , our data annotation protocol aims to automatically obtain 4D annotations of depth map sequence  $\mathcal{D}$  and normal map sequence  $\mathcal{N}$ . Video depth/normal requires consistent predictions between consecutive frames [4, 43]. Therefore, we utilize flow as guidance [39]. By predicting the optical flow, we can guarantee that static areas of a scene remain temporally consistent between consecutive frames. Note that the optical flow method is effective in this case, as existing manipulation datasets [5, 6, 18, 27, 35, 36, 46, 57] are based on fixed camera poses.

To be specific, we employ RAFT [55], an efficient optical flow estimation model, to generate the optical flow  $\mathcal{F}$ . The optical flow between consecutive frames is computed as  $\mathcal{F} = \text{RAFT}(\mathcal{V})$ . The next step involves generating the 3D annotations: depth  $\mathcal{D}$  and surface normal  $\mathcal{N}$ . These 3D annotations provide richer visual context and enhance the understanding of the 4D environment represented by  $\mathcal{V}$ . Recent advances in monocular depth and normal estimators, typically based on models like Vision Transformers (ViT) or Diffusion Models, have shown strong generalization capabilities across diverse datasets.

Then, we leverage the diffusion model Marigold [30] to estimate  $\mathcal{D}^i$  and  $\mathcal{N}^i$  from each video frame  $\mathcal{V}^i$ . The frame is first encoded into a latent representation  $\mathbf{v}^i$  using a variational autoencoder (VAE) [33, 51]:  $\mathbf{v}^i = \text{Encoder}(\mathcal{V}^i)$ . At the same time, an initial depth or normal latent map  $\mathbf{z}_T^i$  is sampled from a Gaussian distribution:  $\mathbf{z}_T^i \sim \mathcal{N}(0, 1)$ . Then we iteratively apply the learned denoiser U-Net [53]  $\epsilon$  on each timestep  $t$  to reconstruct the  $\mathbf{z}_t^i$  and  $\mathcal{Z}^i$ :

$$\mathbf{z}_{t-1}^i = \epsilon(\mathbf{v}^i, \mathbf{z}_t^i, t) \quad \text{and} \quad \mathcal{Z}^i = \text{Decoder}(\mathbf{z}_0^i) \quad (1)$$

where  $(\mathbf{z}, \mathcal{Z}) \in \{(\mathbf{d}, \mathcal{D}), (\mathbf{n}, \mathcal{N})\}$ .

However, most existing approaches [3, 30, 63] are primarily designed for image-based inputs and often struggle when processing videos, resulting in unstable predictions. To overcome this limitation, we propose a novel technique that leverages optical flow as guidance to refine depth estimation. The key insight is that optical flow can capture how static backgrounds and objects move within a scene. Therefore, during the inference stage of the diffusion models, we introduce a loss function that enforces consistency between the backgrounds of consecutive frames. Specifically, we define a static region mask for the  $i$ -th frame based on the optical flow magnitude:  $\mathcal{M}^i = (\|\mathcal{F}^i\| < c)$  where  $c$  is a predefined threshold. In practice, we erode the mask to ensure it covers a robust region. Finally, we define the loss function and integrate its gradient into Eq. 1 for every denoising step:

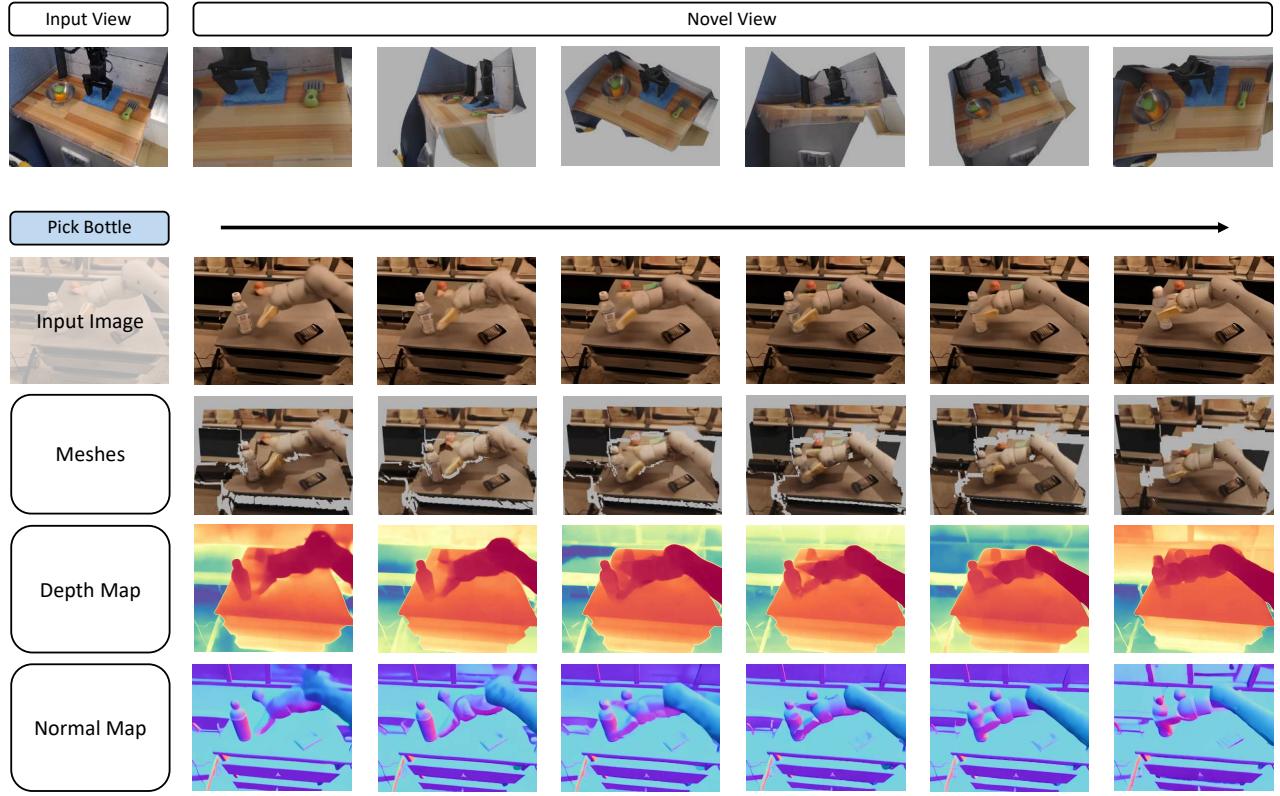
$$\begin{aligned} \mathcal{L}(\mathbf{z}_t^i) = & \left\| \text{Decoder}(\epsilon(\mathbf{v}^i, \mathbf{z}_t^i, t)) \circ (\mathcal{M}^i \cap \mathcal{M}^{i-1}) \right. \\ & \left. - \mathcal{Z}^{i-1} \circ (\mathcal{M}^i \cap \mathcal{M}^{i-1}) \right\|^2 \end{aligned} \quad (2)$$

$$\mathbf{z}_{t-1}^i = \epsilon_\theta \left[ \mathbf{v}^i, \left( \mathbf{z}_t^i - w \nabla_{\mathbf{z}_t^i} \mathcal{L} \right), t \right] \quad (3)$$

where  $w$  is a guidance weight,  $\circ$  represents the element-wise product and  $\mathcal{M}_i \cap \mathcal{M}_{i-1}$  selects the overlapping stable regions between two consecutive frames. The pseudocode is presented in the appendix.

### 3.2. Preliminaries on Video Diffusion Models

Diffusion models [24, 52] are capable of learning the data distribution  $p(x)$  by progressively adding noise to the data until it resembles a Gaussian distribution through a forward process. During inference, a denoiser  $\epsilon$  is trained to recover the data from this noisy state. Latent video diffusion models [70] utilize a Variational Autoencoder (VAE) [33, 56], in the latent space of the data, maintaining high-quality outputs while more efficiently modeling the data distribution. In this section, we formulate the task of RGB  $\mathcal{V}$ , depth  $\mathcal{D}$ , and normal  $\mathcal{N}$  map video generation as a conditional denoising generation task, i.e., we model the distribution  $p(\mathbf{v}, \mathbf{d}, \mathbf{n} | \mathbf{v}_0, \mathcal{T})$ , where  $\mathbf{v}, \mathbf{d}, \mathbf{n}$  represent the predicted future latent sequences of RGB, depth, and normal maps, respectively. The condition  $\mathbf{v}_0$  is a given RGB image latent, and  $\mathcal{T}$  denotes the instruction provided by the user.



**Figure 3. 4D Reconstruction and Video Generation.** The first row shows the results of our 4D reconstruction method from multiple views on a single frame. The following rows present the inference results of our method given the input image and instruction, along with the reconstructed meshes.

225 The forward diffusion process adds Gaussian noise to the  
226 latent  $\mathbf{z} \in \{\mathbf{v}, \mathbf{d}, \mathbf{n}\}$  over  $T$  timesteps, defined as:

$$227 q(\mathbf{z}_t | \mathbf{z}_{t-1}) = \mathcal{N}(\mathbf{z}_t; \sqrt{\alpha_t} \mathbf{z}_{t-1}, (1 - \alpha_t) \mathbf{I}) \quad (4)$$

228 where  $t \in \{1, 2, \dots, T\}$  denotes the diffusion step,  $\alpha_t$  is a  
229 parameter controlling the noise influence at each step, and  $\mathbf{I}$   
230 is the identity matrix. In the reverse process, the model aims  
231 to recover the original latent from the noise. A denoising net-  
232 work  $\epsilon_\theta(\mathbf{x}_t, t, \mathbf{v}_0, \mathcal{T})$  with learning parameters  $\theta$  is trained  
233 to predict the noise added at each timestep. For simplic-  
234 ity, let  $\mathbf{x}_t = [\mathbf{v}_t, \mathbf{n}_t, \mathbf{d}_t]$ , which denotes the concatenation  
235 operation. The reverse process is defined as:

$$236 p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{v}_0, \mathcal{T}) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t, \mathbf{v}_0, \mathcal{T}), \Sigma_\theta(\mathbf{x}_t, t)) \quad (5)$$

237 where

$$238 \mu_\theta(\mathbf{x}_t, t, \mathbf{v}_0, \mathcal{T}) = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \epsilon_\theta(\mathbf{x}_t, t, \mathbf{v}_0, \mathcal{T}) \right). \quad (6)$$

239 The variance term  $\Sigma_\theta(\mathbf{x}_t, t)$  is typically constant. Once  
240 the denoised sequence of latent  $\mathbf{z}_0$  is obtained, the model  
241 reconstructs the final video frames using the decoder net-  
242 work, mapping the latent back to the pixel space:  $\mathcal{Z} =$   
243  $\text{Decoder}(\mathbf{z}_0)$ .

During training, we randomly select samples from the dataset  $(\mathcal{V}, \mathcal{D}, \mathcal{N}, \mathcal{T})$  and apply Eq.4 to add noise  $\epsilon_v$ ,  $\epsilon_d$ , and  $\epsilon_n$  to the RGB-DN data at timestep  $t$ , minimizing the following objective:

$$244 L = \mathbb{E}_{\mathbf{v}_0, \mathcal{T}, t, \epsilon} \left[ \left\| [\epsilon_v, \epsilon_d, \epsilon_n] - \epsilon_\theta(\mathbf{x}_t, t, \mathbf{v}_0, \mathcal{T}) \right\|^2 \right] \quad (7)$$

### 3.3. RGB, Depth and Normal Video Predictions

Training a diffusion model to model temporal RGB-DN data is a challenging task. To effectively train RGB video models, large-scale video datasets with billions of high-quality samples are used [70]. In contrast, even through automatic annotation, our dataset of RGB-DN data contains only around one million data points, which is insufficient to train a world model from scratch. To address this, we finetune the Open-Sora [70] as our RGB-DN prediction model and directly leverage the pre-trained knowledge inside the model to effectively bootstrap our 4D model.

To implement this, we use the temporal VAE [33, 56] in Open-Sora [70] to separately encode RGB, depth, and normal images for each frame of a video, without additional fine-tuning of the VAE. We then expand Open-Sora STDiT's

264 input and output channels threefold. Specifically, we adjust  
 265 the first 3D convolutional input layer and the output linear  
 266 layer to process concatenated RGB-DN video channels. We  
 267 fine-tune all STDiT parameters to denoise videos, allowing  
 268 the model to leverage pre-trained knowledge from video data  
 269 to predict RGB, depth, and normal frames, capturing the 4D  
 270 dynamics of a scene.

### 271 3.4. 4D Scene Reconstruction from RGB-DN Video

272 After obtaining the RGB-DN video, we further optimize  
 273 the depth and reconstruct the surface to generate a full final  
 274 dynamic mesh of the scene. Similar to prior works [30,  
 275 66], our depth representation for each image is given by  
 276 a relative map in the range  $[0, 1]$ , and thus cannot directly  
 277 reconstruct the entire scene. While past work has sidestepped  
 278 this by assuming either a default scale for depth or by directly  
 279 predicting metric depth, such reconstructions from depth are  
 280 often coarse and often cause reconstructed planes or walls to  
 281 be tilted.

282 We instead leverage the normal maps  $\mathcal{N}^i$  and optical  
 283 flow  $\mathcal{F}^i$  between frames in a video to obtain precise per-  
 284 pixel depth estimates from relative depth maps  $\mathcal{D}^i$ . Normal  
 285 maps provide essential information about surface orientation,  
 286 which is vital for enforcing geometric constraints and  
 287 imposing surface smoothness and continuity during depth  
 288 optimization. This spatial optimization leads to more accu-  
 289 rate and reliable depth estimates that closely align with the  
 290 true 3D geometry and capture fine surface details. Addi-  
 291 tionally, we use optical flow between frames to enforce 3D  
 292 consistency over time, effectively optimizing in the temporal  
 293 domain. Together, normals and optical flow offer spatial and  
 294 temporal constraints that enhance depth prediction. Both  
 295 maps, in combination with the coarse depth  $\mathcal{D}^i$ , allow us  
 296 to optimize a refined depth map  $\tilde{\mathcal{D}}$  that corresponds to a  
 297 consistent 4D scene.

298 To formalize the process and enforce consistency across  
 299 frames, we can use the perspective camera model to set  
 300 constraints on the depth and surface normal. In the coordi-  
 301 nate system of the 2D image at frame  $i$ , a pixel position  
 302 is given as  $\mathbf{u} = (u, v)^T \in \mathcal{V}^i$ , and its corresponding depth  
 303 scalar, normal vector is  $d \in \mathcal{D}^i$ ,  $\mathbf{n} = (n_x, n_y, n_z) \in \mathcal{N}^i$ .  
 304 Under the assumption of a perspective camera whose focal  
 305 length is  $f$  and the principal point is  $(c_u, c_v)^T$ , as proposed  
 306 by [15], the log-depth  $\tilde{d} = \log(d)$  should satisfy the fol-  
 307 lowing equations:  $\tilde{n}_z \partial_u \tilde{d} + n_x = 0$  and  $\tilde{n}_z \partial_v \tilde{d} + n_y = 0$   
 308 where  $\tilde{n}_z = n_x(u - c_u) + n_y(v - c_v) + n_z f$ . In addition,  
 309 we can add, the assumption that assumes all locations are  
 310 smooth surfaces [9], we can convert the above constraint to  
 311 the quadratic loss function, allowing us to find the minimized  
 312 depth map:

$$313 \min_{\tilde{d}} \iint_{\Omega} (\tilde{n}_z \partial_u \tilde{d} + n_x)^2 + (\tilde{n}_z \partial_v \tilde{d} + n_y)^2 dudv. \quad (8)$$

314 Following [9], we can convert the above objective to an

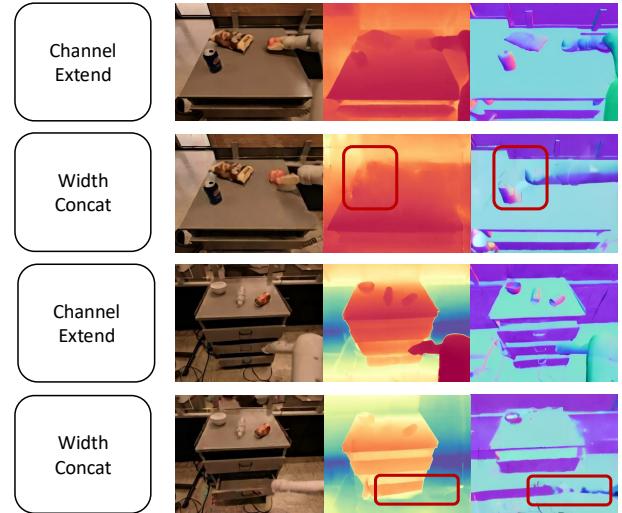


Figure 4. Channel concatenation improves visual quality and ensures consistency across maps.

iteratively optimized loss objective. At iteration step  $t$ , we  
 315 can compute the matrix  $W(\tilde{d}_t)$  and iteratively optimize for  
 316 a refined depth prediction  $\tilde{d}_{t+1}$ :  
 317

$$318 \tilde{d}_{t+1} = \arg \min_{\tilde{d}} (A\tilde{d} - b)^T W(\tilde{d}_t) (A\tilde{d} - b) \stackrel{\text{def}}{=} \arg \min_{\tilde{d}} \mathcal{L}(\tilde{\mathcal{D}}, \mathcal{N}^i), \quad (9)$$

where  $A$  and  $b$  are defined by predicted normals and camera  
 319 intrinsics.  
 320

321 The above approach optimizes depth frame by frame,  
 322 which lacks temporal consistency across the dynamic scene.  
 323 To address this, we compute optical flow between frames  
 324 [55]  $\mathcal{F} = \text{RAFT}(\mathcal{V})$  and enforce consistency of depth across  
 325 frames. We define the static regions of each frame as the pix-  
 326 els with the magnitude of optical flow smaller than threshold  
 327  $\mathcal{M}^i = \|\mathcal{F}^i\| \leq c$ . We then define the dynamic parts of an  
 328 image  $M_d^i$  as all regions not in  $\mathcal{M}^i$ . We further define the  
 329 background of an image  $M_b^i$  as static regions that are fixed  
 330 across image frames,  $M_b^i = \mathcal{M}^i \cap \mathcal{M}^{i-1}$

331 Since optical flow represents the movement of objects in  
 332 the 2D-pixel space, we can retrieve the depth at any position  
 333 from the previous frame to impose consistency constraints.  
 334 To compute the depth values from the previous frame at  
 335 positions corresponding to the current frame, we utilize the  
 336 optical flow  $\mathcal{F}^{i \rightarrow (i-1)}$ . For each pixel  $(u, v)$  in frame  $i$ , the  
 337 optical flow provides the displacement  $(\Delta u, \Delta v)$ , allowing  
 338 us to find the corresponding pixel in frame  $i-1$  at position  
 339  $(u - \Delta u, v - \Delta v)$ . Based on this mapping, we define the  
 340  $\mathcal{D}^{i \rightarrow (i-1)}$  such that:  $\mathcal{D}^{i \rightarrow (i-1)}(u, v) = \mathcal{D}^{i-1}(u - \Delta u, v - \Delta v)$ . We then introduce the loss function  $\mathcal{L}_d$  for dynamic  
 341 regions of an image:  
 342

$$343 \mathcal{L}_d(\tilde{\mathcal{D}}, \hat{\mathcal{D}}^{i-1}, M_d^i, \mathcal{F}^i, \mathcal{F}^{i-1}) = \left\| \tilde{\mathcal{D}}^i \circ M_d^i - \mathcal{D}^{i \rightarrow (i-1)} \circ M_d^i \right\|^2. \quad (10)$$

344 In addition to the loss terms  $\mathcal{L}$  defined previously, we

Method	Chamfer $L_1$	
	RLBench	RT-1
RGB to depth	0.2570	0.3013
4D Point Cloud Diffusion	0.1086	0.2211
Ours	<b>0.0945</b>	<b>0.2022</b>

Table 1. Comparison for 4D Generation on RLBench and RT-1

345 incorporate regularization loss  $\mathcal{L}_g$  enforcing that optimized  
 346 depths are similar to the generated depth map  $\mathcal{D}^i$ . We define  
 347 the regularization loss  $\mathcal{L}_g$  as:

348 
$$\mathcal{L}_g(\mathcal{D}_1, \mathcal{D}_2, \mathcal{M}) = \|\mathcal{D}_1 \circ \mathcal{M} - \mathcal{D}_2 \circ \mathcal{M}\|^2 \quad (11)$$

349 We then define a regularization term on optimized  
 350 depth maps over background regions of images  
 351  $\lambda_{g2}\mathcal{L}_g(\tilde{\mathcal{D}}, \hat{\mathcal{D}}^{i-1}, \mathcal{M}_b^i)$  enforcing that optimized depths for  
 352 background regions of an image are consistent between  
 353 frames and a regularization term over dynamic regions of  
 354 images  $\lambda_{g1}\mathcal{L}_g(\tilde{\mathcal{D}}, \mathcal{D}^i, M_d^i)$  enforcing that optimized depth  
 355 of dynamic regions of an image match with predicted  
 356 dynamic depths.

357 The overall loss objective we optimize is given by:

358 
$$\arg \min_{\tilde{\mathcal{D}}} \mathcal{L}(\tilde{\mathcal{D}}, \mathcal{N}^i) + \lambda_b \mathcal{L}_g(\tilde{\mathcal{D}}, \hat{\mathcal{D}}^{i-1}, \mathcal{M}_b^i) + \quad (12)$$

359 
$$+ \lambda_d \mathcal{L}_d(\tilde{\mathcal{D}}, \hat{\mathcal{D}}^{i-1}, \mathcal{M}_d^i, \mathcal{F}^i, \mathcal{F}^{i-1}) + \quad (13)$$

360 
$$+ \lambda_{g2} \mathcal{L}_g(\tilde{\mathcal{D}}, \mathcal{D}^i, M_d^i) + \lambda_{g1} \mathcal{L}_g(\tilde{\mathcal{D}}, \mathcal{D}^i, M_d^i). \quad (14)$$

361 We initialize the starting depth  $\tilde{d}_0 = \mathcal{D}^i$  with the generated  
 362 depth map, and similar to prior work [9, 62], where we  
 363 repeatedly optimize  $\tilde{\mathcal{D}}^t$  across multiple iterations (using the  
 364 previously optimized depth map  $\tilde{\mathcal{D}}^{t-1}$  to define the new  
 365 optimization objective).

366 Finally, we construct faces by connecting pixels to their  
 367 nearby neighbors, resulting in a mesh derived from the op-  
 368 timized depth. Meshes provide a structured representation  
 369 of the 3D geometry of a scene, enabling detailed surface  
 370 reconstruction and facilitating further processing like ren-  
 371 dering or physical simulations. For mesh denoising, we  
 372 remove isolated vertices based on mean neighbor distances  
 373 and eliminate small clusters using DBSCAN [22]. We also  
 374 discard faces with abnormal normals or high edge-length  
 375 variance, ensuring the final mesh is cleaner and more suitable  
 376 for downstream tasks.

### 377 3.5. Inverse Dynamics Models from 4D Scenes

378 After generating 4D scenes, which encapsulate both spatial  
 379 and temporal information, we extract geometric details that  
 380 can significantly enhance downstream tasks in robotics. The  
 381 detailed geometry captured by these scenes plays a crucial  
 382 role in robotic grasping tasks.

Method	Time Cost	Consistency	
		Depth	Normal
Open-Sora	25.6 seconds	0.09267	0.04153
Marigold-LCM	<b>15.2 seconds</b>	0.09453	0.04647
Guided Marigold	~3.5 hours	<b>0.07299</b>	<b>0.03822</b>

Table 2. Comparison of Data Generation Methods

To achieve this, we employ an inverse dynamics model built on the 4D meshes, predicting the appropriate robot action  $a_i$  based on the current state  $s_i$ , the predicted future state  $s_{i+1}$  and the instruction  $\mathcal{T}$ . Mathematically, this relationship is expressed as  $a_i = \text{ID}(s_i, s_{i+1}, \mathcal{T})$ . In our scenario,  $s_i$  represents the scene at time step  $i$ . Specifically, we sample the meshes to obtain point clouds, which are encoded by a PointNet [48] architecture within the inverse dynamics model to extract features. These features, combined with the instruction text embeddings, are further processed by an MLP to generate the final action.

## 4. Experiments

In this section, we evaluate the performance of our proposed model across several tasks. In Section 4.1, we present our experiments on 4D mesh prediction using the RLBench [27] and RT-1 [5] datasets. In Section 4.2, we conduct experiments on embodied novel view synthesis, using RLBench to assess our model’s ability to generate novel views from monocular video inputs. In Section 4.3, we explore embodied action planning, applying our model to guide robotic arm policies for specific tasks. Finally, in Section 4.4, we present discussions and ablation studies that analyze the effect of different architectural and data generation choices on the quality and consistency of our video diffusion models. For additional experiments and results, see the Appendix.

### 4.1. 4D Scene Prediction

Since no prior work directly generates dynamic meshes from the first frame image and text inputs, we primarily compare our method to a 4D point cloud diffusion model. **Baseline:** our baselines include two main approaches. (1) The first is a 16-frame RGB diffusion model, where we obtain depth from the pre-trained depth estimator Marigold [30], and lift it to 3D via camera intrinsic and extrinsic parameters. (2) We also modify the Point-E [45] model by conditioning it on the mean of CLIP [49] features extracted from both text and image inputs, outputting a point cloud of size  $T \times \text{num of points}$ , where  $T$  is set to 4 due to computational constraints. **Dataset:** the datasets used for evaluation include RLBench [27] and our annotated RT-1 [5] dataset. **Metric:** for evaluation, we use the  $L_1$  Chamfer Distance metric, which measures the distance between two point sets. The results are shown in Table 1.

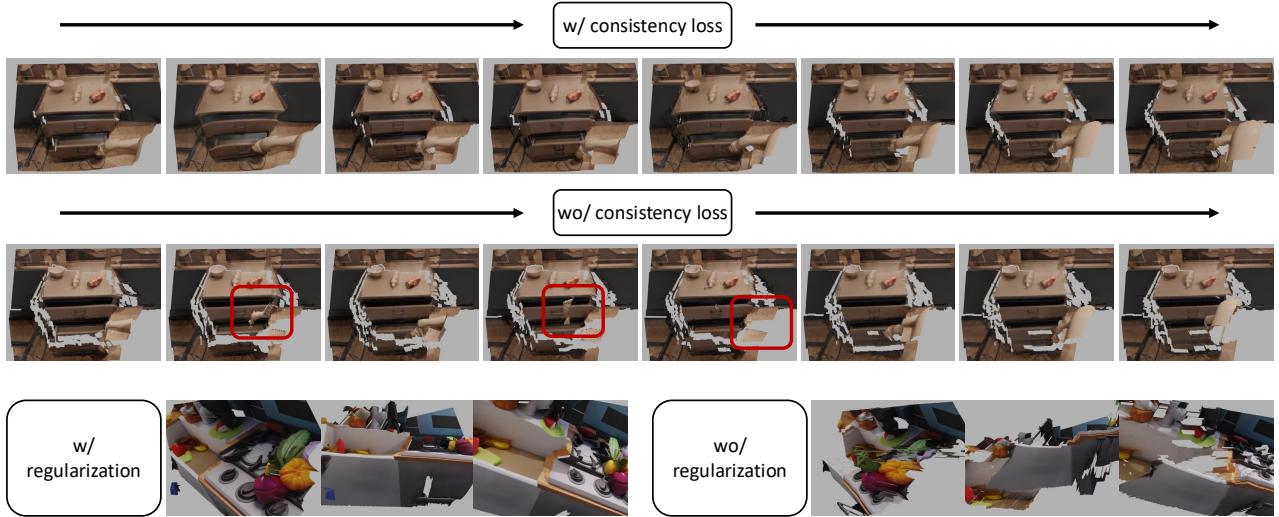


Figure 5. Effect of Consistency and Regularization Losses on 4D Mesh Reconstruction. The red boxes highlight the inconsistent regions.

As shown in the table, our method achieves the lowest Chamfer distances on both the RLBench [27] and RT-1 [5] datasets, indicating a more accurate reconstruction of 4D structures compared to the baselines. All methods perform better on RLBench [27], which is due to it being synthetic data, with less noise and perfectly accurate depth ground truth. The RGB-to-depth approach, while simple, suffers from larger errors due to the limitations of depth estimation from 2D images. The 4D point cloud diffusion method performs better, particularly on RLBench, but still lags behind our approach. Additionally, point cloud training is computationally expensive, restricting the number of frames used. In contrast, our model, by leveraging both image and text inputs, manages to generate more precise 3D representations, particularly in capturing fine-grained details in dynamic scenes. We show our qualitative results in Figure 3.

## 4.2. Embodied Novel View Synthesis

Our method performs monocular video to 4D tasks by predicting depth and normal sequences and generating meshes. **Baseline:** (1) we select Shape of Motion [58] as our primary baseline, a state-of-the-art video reconstruction approach that utilizes Gaussian splatting [31]. (2) Additionally, we include an RGB-to-depth approach, lifting depth to 3D and rendering point clouds for novel views. **Dataset:** since real-world datasets like RT-1 lack multiview camera information, we conduct experiments on RLBench. The input is a monocular front camera video, and we compare results from the overhead and left shoulder cameras. **Metric:** PSNR (reconstruction accuracy), SSIM (structural similarity), LPIPS (perceptual difference), CLIP Score (semantic match) [71], CLIP aesthetic (visual quality) [37], and Time costs.

Method	PSNR	SSIM	LPIPS	CLIP Score	CLIP aesthetic	Time costs
Shape of Motion	10.94	0.2402	<b>0.7382</b>	66.67	3.61	~2 hours
Ours	<b>12.99</b>	<b>0.4262</b>	0.6051	<b>83.02</b>	<b>3.73</b>	~ 1 minutes

Table 3. Performance Comparison of Novel View Synthesis Methods on RLBench Dataset

## 4.3. Embodied Action Planning

Since our world model can predict future scenes, a direct application is to guide robotic arm policies. **Baseline:** we compare our method with re-implemented UniPi\*, which is a 2D world model and uses an inverse dynamic policy to predict actions based on the videos. For this baseline, we fine-tune OpenSora [70] on RLBench data as the 2D world model and use a ResNet [23] in the image-based inverse dynamic model to encode both the current and predicted frames. Then, the encoded features and the text embeddings are passed through an MLP to obtain the 7-DoF action. **Dataset:** for simulation, we use RLBench [27] and collect 500 samples for each task to train our model. **Metric:** the success rate, which is returned by the simulator. For both the baseline and our method, we train the inverse dynamic model on collected RLBench data, with noise added to the corresponding modalities (image/3D point cloud). Given an initial state during inference, we first predict and record all future keyframes. In subsequent actions, we only query the inverse dynamic model to obtain the corresponding actions by the current state and the predicted future state. Table 5 reports the accuracy across different tasks.

The results show that our method outperforms video diffusion models across the selected tasks. This is because, in most tasks, 4D meshes or point clouds can reveal the geometry of objects, providing better spatial guidance for robotics planning, as seen in tasks like Close Box and Close Laptop. However, the performance of our model declines in tasks

Frames	Channel	RGB			AbsRel ↓	Depth		Normal		Consistency Edge-sim ↑
		FVD ↑	SSIM ↑	PSNR ↑		$\delta_1$ ↑	$\delta_2$ ↑	Mean ↓	Median ↓	
32	✗	20.12	70.23	19.32	30.22	59.21	80.28	54.22	40.87	6.41
32	✓	19.84	69.94	19.30	18.67	69.65	89.12	19.78	10.01	26.42
16	✗	<b>25.78</b>	71.89	21.86	<b>16.14</b>	76.59	91.54	26.59	15.73	24.36
16	✓	25.45	<b>74.98</b>	<b>21.94</b>	16.53	<b>77.13</b>	<b>92.15</b>	<b>16.23</b>	<b>7.78</b>	<b>38.50</b>
										<b>32.98</b>

Table 4. Impact of the number of frames and channel concatenation. Frames refer to the number of input frames used in the model. “Channel” ✓ indicates concatenating RGB, depth, and normal maps along the channel dimension, while ✗ refers to concatenation along the width.

Methods	Close Box	Close Laptop	Lift Block	Lamp Off
UniPi*	81.0	14.3	19.0	<b>57.1</b>
Ours	<b>95.2</b>	<b>28.6</b>	<b>23.8</b>	42.9

Table 5. Evaluation of action planning on RLBench dataset

like Lamp Off, due to the small size of the switch, which may not have been sampled. Overall, these results highlight the potential of combining 4D scene prediction with inverse dynamic models to improve robotics task execution.

#### 4.4. Ablations

**RGB-DN Video Diffusion Models** We first conduct ablation studies on our video generation task. We perform experiments to explore the impact of different concatenation methods and the number of frames on the results. For the former, we compare two settings: (1) concatenating RGB-DN images along the width to form a larger image, or (2) using a VAE encoder [34, 56] to separately process RGB, depth, and normal maps, and concatenating them along the channel dimension before inputting them into the diffusion model. In the latter case, we also modify the input and output dimensions of the backbone network. Our evaluation metrics focus on the generation/reconstruction quality of RGB, depth, and normal maps. Additionally, we introduce an edge similarity metric to assess the consistency across RGB, depth, and normal maps at the same timestamp. Specifically, we convert them into gray-scale images, apply Canny edge detection [8], and compare the edge maps using SSIM [59].

Although concatenating images along the width results in better RGB reconstruction due to better utilization of the pre-trained Open-Sora model [70], it is less effective for depth and normal map predictions. Moreover, the inconsistency between RGB, depth, and normal maps prevents effective post-processing. As shown in Figure 4, concatenating along the channel dimension yields higher-quality depth and normal maps while maintaining consistency across the three value maps. This prevents issues such as a robotic arm appearing in different locations in the RGB and depth/normal maps.

**Data Collection** This part primarily compares the effects of guidance on depth and normal diffusion models during data generation. As shown in Table 2, “Marigold-LCM” refers to our use of the Marigold Latent Consistency Model

[30] for independently predicting each frame. “Guided Marigold” represents our data generation method. We compare the time cost and the static part  $L_1$  difference for these methods. As a reference, we provide the scores from our trained Open-Sora. Qualitative results are presented in the Appendix, where we observe that our proposed data generation method maintains the highest consistency, though at a significantly higher time cost. This also highlights the necessity of training a world model to rapidly predict and generate dynamic scenes.

**Regularization and Consistency Loss** in 4D Mesh Reconstruction. In this task, we evaluate the impact of our newly designed loss terms, as shown in Figure 5. The first two rows demonstrate the effect of the consistency loss, where we render frames from the same camera view at different time steps. The results show that the robot arm’s movements are more coherent with the consistency loss applied. The last row highlights the role of the regularization loss. We display images of the same frame from three different views, revealing that this loss term helps improve the geometric accuracy of the reconstruction.

#### 5. Conclusion

Our current approach has several limitations. First, while our RGB-DN representation of a 4D world model is cheap and easy to predict, it only captures a single surface of the world. To construct a more complete 4D world model, it may be interesting in the future to have a generative model that generates multiple RGB-DN views of the world, which can then be integrated to form a more complete 4D world model. In addition, we observe that generated RGB-DN maps from our 4D world model may not be fully consistent with each. Adding additional structure in the architecture or loss of function constraints at training time to help enforce consistency is a rich direction for future work.

Overall, our work provides some first steps towards the goal of constructing a 4D generative model of the world. We believe that such world models will be increasingly powerful and useful in the future, serving as a way to simulate the physical world and is an important step towards constructing intelligent embodied agents. Such models would then enable us to train policies in the real world in a fully offline manner, as well as roll out and imagine future plans in the world.

563 **References**

- [1] Alessandro Achille and Stefano Soatto. A separation principle for control in the age of deep learning. *Annual Review of Control, Robotics, and Autonomous Systems*, 1:287–307, 2018. 2
- [2] Dimitri P. Bertsekas. *Dynamic Programming and Optimal Control*. Athena Scientific, 1995. 2
- [3] Shariq Farooq Bhat, Reiner Birkl, Diana Wofk, Peter Wonka, and Matthias Müller. Zoedepth: Zero-shot transfer by combining relative and metric depth. *arXiv preprint arXiv:2302.12288*, 2023. 3
- [4] Jiawang Bian, Zhichao Li, Naiyan Wang, Huangying Zhan, Chunhua Shen, Ming-Ming Cheng, and Ian Reid. Unsupervised scale-consistent depth and ego-motion learning from monocular video. *Advances in neural information processing systems*, 32, 2019. 3
- [5] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022. 3, 6, 7
- [6] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*, 2023. 2, 3
- [7] Jake Bruce, Michael D Dennis, Ashley Edwards, Jack Parker-Holder, Yuge Shi, Edward Hughes, Matthew Lai, Aditi Mavalankar, Richie Steigerwald, Chris Apps, et al. Genie: Generative interactive environments. In *Forty-first International Conference on Machine Learning*, 2024. 2
- [8] John Canny. A computational approach to edge detection. *IEEE Transactions on pattern analysis and machine intelligence*, (6):679–698, 1986. 8
- [9] Xu Cao, Hiroaki Santo, Boxin Shi, Fumio Okura, and Yasuyuki Matsushita. Bilateral normal integration. In *European Conference on Computer Vision*, pages 552–567. Springer, 2022. 5, 6
- [10] Chang Chen, Yi-Fu Wu, Jaesik Yoon, and Sungjin Ahn. Transdreamer: Reinforcement learning with transformer world models. *arXiv preprint arXiv:2202.09481*, 2022. 2
- [11] Silvia Chiappa, Sébastien Racaniere, Daan Wierstra, and Shakir Mohamed. Recurrent environment simulators. *arXiv preprint arXiv:1704.02254*, 2017. 2
- [12] Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Akanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. Palm-e: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023. 2
- [13] Yilun Du, Mengjiao Yang, Pete Florence, Fei Xia, Ayzaan Wahid, Brian Ichter, Pierre Sermanet, Tianhe Yu, Pieter Abbeel, Joshua B Tenenbaum, et al. Video language planning. *arXiv preprint arXiv:2310.10625*, 2023. 1, 2
- [14] Yilun Du, Sherry Yang, Bo Dai, Hanjun Dai, Ofir Nachum, Josh Tenenbaum, Dale Schuurmans, and Pieter Abbeel. Learning universal policies via text-guided video generation. *Advances in Neural Information Processing Systems*, 36, 2024. 1, 2
- [15] Jean-Denis Durou, Jean-François Aujol, and Frédéric Courteille. Integrating the normal field of a surface in the presence of discontinuities. In *International Workshop on Energy Minimization Methods in Computer Vision and Pattern Recognition*, pages 261–273. Springer, 2009. 5
- [16] Norm Ferns, Prakash Panangaden, and Doina Precup. Metrics for finite markov decision processes. In *UAI*, pages 162–169, 2004. 2
- [17] Roya Firoozi, Johnathan Tucker, Stephen Tian, Anirudha Majumdar, Jiankai Sun, Weiyu Liu, Yuke Zhu, Shuran Song, Ashish Kapoor, Karol Hausman, et al. Foundation models in robotics: Applications, challenges, and the future. *The International Journal of Robotics Research*, page 02783649241281508, 2023. 2
- [18] Susan R Fussell, Leslie D Setlock, and Robert E Kraut. Effects of head-mounted and scene-oriented video systems on remote collaboration on physical tasks. In *Proceedings of the SIGCHI conference on Human factors in computing systems*, pages 513–520, 2003. 3
- [19] Maitrey Gramopadhye and Daniel Szafrir. Generating executable action plans with environmentally-aware language models. *arXiv preprint arXiv:2210.04964*, 2022. 2
- [20] David Ha and Jürgen Schmidhuber. Recurrent world models facilitate policy evolution. In *Advances in Neural Information Processing Systems*. Curran Associates, Inc., 2018. 1, 2
- [21] Danijar Hafner, Timothy P Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering atari with discrete world models. In *International Conference on Learning Representations*, 2021. 2
- [22] Michael Hahsler, Matthew Piekenbrock, and Derek Doran. dbscan: Fast density-based clustering with R. *Journal of Statistical Software*, 91(1):1–30, 2019. 6
- [23] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 7
- [24] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 3
- [25] Yining Hong, Haoyu Zhen, Peihao Chen, Shuhong Zheng, Yilun Du, Zhenfang Chen, and Chuang Gan. 3d-lm: Injecting the 3d world into large language models. *Advances in Neural Information Processing Systems*, 36:20482–20494, 2023. 2
- [26] Jiangyong Huang, Silong Yong, Xiaoqian Ma, Xiongkun Linghu, Puha Li, Yan Wang, Qing Li, Song-Chun Zhu, Baoxiong Jia, and Siyuan Huang. An embodied generalist agent in 3d world. *arXiv preprint arXiv:2311.12871*, 2023. 2
- [27] Stephen James, Zicong Ma, David Rovick Arrojo, and Andrew J. Davison. Rlbench: The robot learning benchmark & learning environment. *IEEE Robotics and Automation Letters*, 2020. 3, 6, 7
- [28] Michael Janner, Yilun Du, Joshua B Tenenbaum, and Sergey Levine. Planning with diffusion for flexible behavior synthesis. *arXiv preprint arXiv:2205.09991*, 2022. 1

- 675 [29] Yunfan Jiang, Agrim Gupta, Zichen Zhang, Guanzhi Wang,  
676 Yongqiang Dou, Yanjun Chen, Li Fei-Fei, Anima Anand-  
677 kumar, Yuke Zhu, and Linxi Fan. Vima: General robot  
678 manipulation with multimodal prompts. *arXiv preprint*  
679 *arXiv:2210.03094*, 2(3):6, 2022. 2
- 680 [30] Bingxin Ke, Anton Obukhov, Shengyu Huang, Nando Met-  
681 zger, Rodrigo Caye Daudt, and Konrad Schindler. Repurposing  
682 diffusion-based image generators for monocular depth  
683 estimation. In *Proceedings of the IEEE/CVF Conference on*  
684 *Computer Vision and Pattern Recognition*, pages 9492–9502,  
685 2024. 2, 3, 5, 6, 8
- 686 [31] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and  
687 George Drettakis. 3d gaussian splatting for real-time radiance  
688 field rendering. *ACM Trans. Graph.*, 42(4):139–1, 2023. 7
- 689 [32] Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted  
690 Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailev, Ethan  
691 Foster, Grace Lam, Pannag Sanketi, et al. Openvla: An  
692 open-source vision-language-action model. *arXiv preprint*  
693 *arXiv:2406.09246*, 2024. 2
- 694 [33] Diederik P Kingma. Auto-encoding variational bayes. *arXiv*  
695 *preprint arXiv:1312.6114*, 2013. 3, 4
- 696 [34] Diederik P Kingma. Auto-encoding variational bayes. *arXiv*  
697 *preprint arXiv:1312.6114*, 2013. 8
- 698 [35] Danica Kragic, Mårten Björkman, Henrik I Christensen, and  
699 Jan-Olof Eklundh. Vision for robotic object manipulation in  
700 domestic settings. *Robotics and autonomous Systems*, 52(1):  
701 85–100, 2005. 3
- 702 [36] Nikolaos Kyriazis and Antonis Argyros. Physically plausible  
703 3d scene tracking: The single actor hypothesis. In *Proceed-  
704 ings of the IEEE conference on computer vision and pattern  
705 recognition*, pages 9–16, 2013. 3
- 706 [37] LAION-AI. Aesthetic predictor. [https://github.com/  
707 LAION-AI/aesthetic-predictor](https://github.com/LAION-AI/aesthetic-predictor), 2022. 7
- 708 [38] Timothée Lesort, Natalia Díaz-Rodríguez, Jean-François  
709 Goudou, and David Filliat. State representation learning  
710 for control: An overview. *Neural Networks*, 108:379–392,  
711 2018. 2
- 712 [39] Siyuan Li, Yue Luo, Ye Zhu, Xun Zhao, Yu Li, and Ying Shan.  
713 Enforcing temporal consistency in video depth estimation. In  
714 *Proceedings of the IEEE/CVF International Conference on*  
715 *Computer Vision*, pages 1145–1154, 2021. 3
- 716 [40] Shuang Li, Xavier Puig, Chris Paxton, Yilun Du, Clinton  
717 Wang, Linxi Fan, Tao Chen, De-An Huang, Ekin Akyürek,  
718 Anima Anandkumar, et al. Pre-trained language models for  
719 interactive decision-making. *Advances in Neural Information  
720 Processing Systems*, 35:31199–31212, 2022. 2
- 721 [41] Junbang Liang, Ruoshi Liu, Ege Ozguroglu, Sruthi Sudhakar,  
722 Achal Dave, Pavel Tokmakov, Shuran Song, and Carl Von-  
723 drick. Dreamitate: Real-world visuomotor policy learning via  
724 video generation. *arXiv preprint arXiv:2406.16862*, 2024. 1,  
725 2
- 726 [42] Lennart Ljung and Torkel Glad. *Modeling of dynamic systems*.  
727 Prentice-Hall, Inc., 1994. 2
- 728 [43] Fei Luo, Lin Wei, and Chunxia Xiao. Stable depth estimation  
729 within consecutive video frames. In *Advances in Computer  
730 Graphics: 38th Computer Graphics International Conference,  
731 CGI 2021, Virtual Event, September 6–10, 2021, Proceedings*  
732 38, pages 54–66. Springer, 2021. 3
- 733 [44] Vincent Micheli, Eloi Alonso, and François Fleuret. Trans-  
734 formers are sample efficient world models. *arXiv preprint*  
735 *arXiv:2209.00588*, 2022. 2
- 736 [45] Alex Nichol, Heewoo Jun, Prafulla Dhariwal, Pamela  
737 Mishkin, and Mark Chen. Point-e: A system for generat-  
738 ing 3d point clouds from complex prompts. *arXiv preprint*  
739 *arXiv:2212.08751*, 2022. 6
- 740 [46] Antje Nuthmann and Teresa Canas-Bajo. Visual search in  
741 naturalistic scenes from foveal to peripheral vision: A compar-  
742 ision between dynamic and static displays. *Journal of Vision*,  
743 22(1):10–10, 2022. 3
- 744 [47] OpenAI. Video generation models as world simu-  
745 lators. [https://openai.com/index/video-  
746 generation-models-as-world-simulators/](https://openai.com/index/video-generation-models-as-world-simulators/), 2023. Accessed: 2024-10-01. 2
- 747 [48] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J  
748 Guibas. Pointnet++: Deep hierarchical feature learning on  
749 point sets in a metric space. *Advances in neural information  
750 processing systems*, 30, 2017. 6
- 751 [49] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya  
752 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,  
753 Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning  
754 transferable visual models from natural language supervi-  
755 sion. In *International conference on machine learning*, pages  
756 8748–8763. PMLR, 2021. 6
- 757 [50] Shreyas Sundara Raman, Vanya Cohen, Eric Rosen, Ifrah  
758 Idrees, David Paulius, and Stefanie Tellex. Planning with  
759 large language models via corrective re-prompting. *arXiv*  
760 *preprint arXiv:2211.09935*, 2022. 2
- 761 [51] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wier-  
762 stra. Stochastic backpropagation and approximate inference  
763 in deep generative models. In *International conference on*  
764 *machine learning*, pages 1278–1286. PMLR, 2014. 3
- 765 [52] Robin Rombach, Andreas Blattmann, Dominik Lorenz,  
766 Patrick Esser, and Björn Ommer. High-resolution image  
767 synthesis with latent diffusion models. In *Proceedings of*  
768 *the IEEE/CVF conference on computer vision and pattern  
769 recognition*, pages 10684–10695, 2022. 3
- 770 [53] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-  
771 net: Convolutional networks for biomedical image segmen-  
772 tation. In *Medical image computing and computer-assisted  
773 intervention—MICCAI 2015: 18th international conference,  
774 Munich, Germany, October 5–9, 2015, proceedings, part III  
775 I8*, pages 234–241. Springer, 2015. 3
- 776 [54] Richard S Sutton. Dyna, an integrated architecture for learn-  
777 ing, planning, and reacting. *ACM Sigart Bulletin*, 2(4):160–  
778 163, 1991. 2
- 779 [55] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field  
780 transforms for optical flow. In *Computer Vision—ECCV 2020:  
781 16th European Conference, Glasgow, UK, August 23–28,  
782 2020, Proceedings, Part II 16*, pages 402–419. Springer, 2020.  
783 3, 5
- 784 [56] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete  
785 representation learning. *Advances in neural information pro-  
786 cessing systems*, 30, 2017. 3, 4, 8
- 787 [57] Homer Walke, Kevin Black, Abraham Lee, Moo Jin Kim,  
788 Max Du, Chongyi Zheng, Tony Zhao, Philippe Hansen-  
789 Estruch, Quan Vuong, Andre He, Vivek Myers, Kuan Fang,  
790 2

- 791 Chelsea Finn, and Sergey Levine. Bridgedata v2: A dataset  
792 for robot learning at scale. In *Conference on Robot Learning*  
793 (*CoRL*), 2023. 3
- 794 [58] Qianqian Wang, Vickie Ye, Hang Gao, Jake Austin, Zhengqi  
795 Li, and Angjoo Kanazawa. Shape of motion: 4d reconstruc-  
796 tion from a single video. *arXiv preprint arXiv:2407.13764*,  
797 2024. 7
- 798 [59] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P  
799 Simoncelli. Image quality assessment: from error visibility to  
800 structural similarity. *IEEE transactions on image processing*,  
801 13(4):600–612, 2004. 8
- 802 [60] Zihao Wang, Shaofei Cai, Guanzhou Chen, Anji Liu, Xiaojian  
803 Ma, and Yitao Liang. Describe, explain, plan and select:  
804 interactive planning with llms enables open-world multi-task  
805 agents. In *Thirty-seventh Conference on Neural Information  
806 Processing Systems*, 2023. 2
- 807 [61] Jiannan Xiang, Guangyi Liu, Yi Gu, Qiyue Gao, Yuting Ning,  
808 Yuheng Zha, Zeyu Feng, Tianhua Tao, Shibo Hao, Yemin  
809 Shi, et al. Pandora: Towards general world model with  
810 natural language actions and video states. *arXiv preprint  
811 arXiv:2406.09455*, 2024. 1, 2
- 812 [62] Yuliang Xiu, Jinlong Yang, Xu Cao, Dimitrios Tzionas, and  
813 Michael J. Black. ECON: Explicit Clothed humans Optimized  
814 via Normal integration. In *Proceedings of the IEEE/CVF Con-  
815 ference on Computer Vision and Pattern Recognition (CVPR)*,  
816 2023. 6
- 817 [63] Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiashi  
818 Feng, and Hengshuang Zhao. Depth anything: Unleashing  
819 the power of large-scale unlabeled data. In *CVPR*, 2024. 2, 3
- 820 [64] Mengjiao Yang, Yilun Du, Kamyar Ghasemipour, Jonathan  
821 Tompson, Dale Schuurmans, and Pieter Abbeel. Learn-  
822 ing interactive real-world simulators. *arXiv preprint  
823 arXiv:2310.06114*, 2023. 1, 2
- 824 [65] Sherry Yang, Ofir Nachum, Yilun Du, Jason Wei, Pieter  
825 Abbeel, and Dale Schuurmans. Foundation models for deci-  
826 sion making: Problems, methods, and opportunities. *arXiv  
827 preprint arXiv:2303.04129*, 2023. 2
- 828 [66] Chongjie Ye, Lingteng Qiu, Xiaodong Gu, Qi Zuo, Yushuang  
829 Wu, Zilong Dong, Liefeng Bo, Yuliang Xiu, and Xiaoguang  
830 Han. Stablenormal: Reducing diffusion variance for stable  
831 and sharp normal. *arXiv preprint arXiv:2406.16864*, 2024. 5
- 832 [67] Hongxin Zhang, Weihua Du, Jiaming Shan, Qinrong Zhou,  
833 Yilun Du, Joshua B Tenenbaum, Tianmin Shu, and Chuang  
834 Gan. Building cooperative embodied agents modularly with  
835 large language models. *arXiv preprint arXiv:2307.02485*,  
836 2023. 2
- 837 [68] Hongxin Zhang, Zeyuan Wang, Qiushi Lyu, Zheyuan Zhang,  
838 Sunli Chen, Tianmin Shu, Yilun Du, and Chuang Gan. Combo:  
839 Compositional world models for embodied multi-  
840 agent cooperation. *arXiv preprint arXiv:2404.10775*, 2024.  
841 1, 2
- 842 [69] Haoyu Zhen, Xiaowen Qiu, Peihao Chen, Jincheng Yang, Xin  
843 Yan, Yilun Du, Yining Hong, and Chuang Gan. 3d-vla: A  
844 3d vision-language-action generative world model. *arXiv  
845 preprint arXiv:2403.09631*, 2024. 2
- 846 [70] Zangwei Zheng, Xiangyu Peng, Tianji Yang, Chenhui Shen,  
847 Shenggui Li, Hongxin Liu, Yukun Zhou, Tianyi Li, and Yang  
You. Open-sora: Democratizing efficient video production  
848 for all, 2024. 1, 2, 3, 4, 7, 8
- [71] SUN Zhengwendai. clip-score: CLIP Score for Py-  
Torch. <https://github.com/taited/clip-score>, 2023. Version 0.1.1. 7
- [72] Siyuan Zhou, Yilun Du, Jiaben Chen, Yandong Li, Dit-Yan  
Yeung, and Chuang Gan. Robodreamer: Learning composi-  
tional world models for robot imagination. *arXiv preprint  
arXiv:2404.12377*, 2024. 2
- [73] Fangqi Zhu, Hongtao Wu, Song Guo, Yuxiao Liu, Chilam  
Cheang, and Tao Kong. Irasim: Learning interactive real-  
robot action simulators. *arXiv preprint arXiv:2406.14540*,  
2024. 1
- [74] Karl J. Åström and Björn Wittenmark. Adaptive control of  
linear time-invariant systems. *Automatica*, 9(6):551–564,  
1973. 2