

Learning 4D Embodied World Models

Anonymous CVPR submission

Paper ID 11899

Abstract

In this paper, we present an effective approach for learning novel 4D embodied world models, which predict the dynamic evolution of 3D scenes over time in response to an embodied agent’s actions, providing both spatial and temporal consistency. We propose to learn a 4D world model by training on RGB-DN (RGB, Depth, and Normal) data. This not only surpasses traditional 2D models by incorporating detailed shape, configuration, and temporal changes into their predictions, but also allows us to effectively learn accurate inverse dynamic models for an embodied agent. Specifically, we first annotate temporally coherent depth and normal information for existing robotic manipulation video datasets leveraging both off-the-shelf models and optical flow guidance techniques. Next, we fine-tune a video generation 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. Our method enables us to predict high-quality meshes coherent 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 released soon.

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).

A key challenge for training a 4D world model is a lack of access to existing large-scale datasets with existing 4D

038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064

065
066
067
068
069
070
071
072
073
074
075
076
077
078

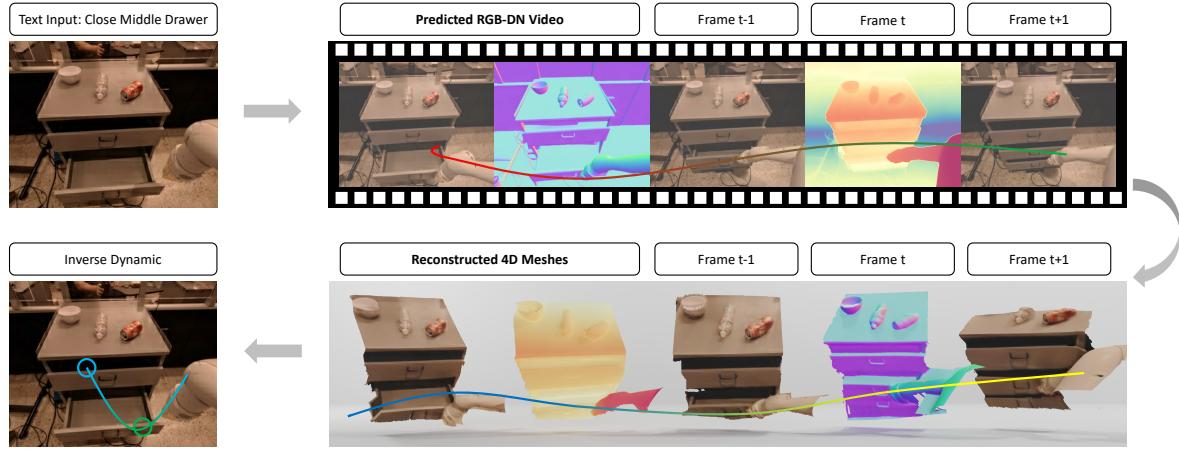


Figure 1. 4D Embodied World Models. Our approach gets an input image and instructions and predicts RGB-D-N (RGB, Depth, Normal) maps. We propose a normal integration method that constructs a high-quality 4D mesh of the interaction from these predictions.

annotations, or the high-quality image, depth, and normal annotations needed to train our approach. To construct such data, in principle, we can use pre-trained depth and normal estimators [30, 63] to obtain such estimates from video data, but these estimators are typically limited to predicting relative value maps from individual frames. As a result, when a scene changes, these estimators will produce depth and normal maps with inconsistencies across time. To tackle this challenge, we develop a data collection pipeline that leverages optical flow between frames in a video to enforce consistency between generated depth and normal maps across timesteps. In particular, we use optical flow to guide a depth and normal diffusion model across frames in a video, which we find is sufficient to ensure consistent depth and normal predictions across all timesteps without the need for expensive ground-truth annotations, facilitating the training of our world model on a large scale.

Overall, our paper has the following contributions: **(1)** We introduce a 4D embodied world model, and present an efficient representation of this model, in the form of RGB, depth, and normal maps, and illustrate how this representation can be used to construct a full 4D scene. **(2)** We present a pipeline to automatically extract 4D world model data from existing robot video datasets, leveraging an optical flow-guided depth and normal diffusion model. **(3)** We illustrate the approach's efficacy in generating consistent 4D meshes across different environments, substantially outperforming other baselines, and showing its downstream use as an effective policy.

2. Related Work

Embodied Foundation Models A flurry of recent work has focused on constructing foundation models for general purpose agents [17, 65]. One line of work has focused on constructing multimodal language models that operate over

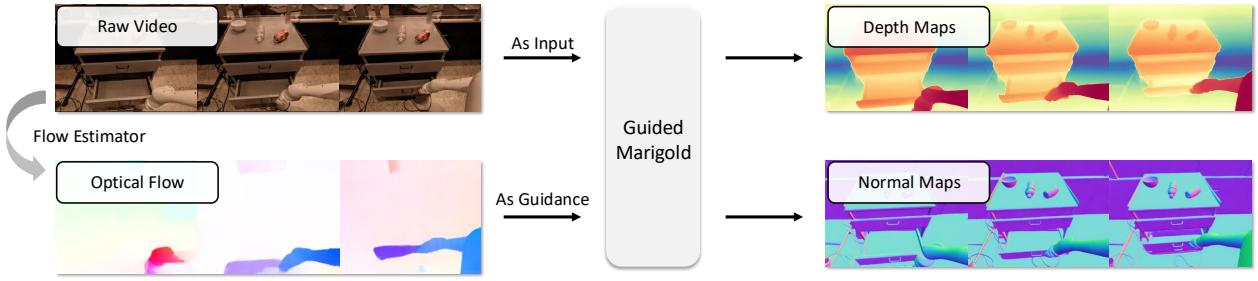
images [12, 19, 29, 40, 50, 60, 67] as well as 3D inputs [25, 26] and output text describing the actions of an agent. Other works have focused on the construction of vision-language-action (VLA) models that directly output action tokens [6, 32, 69]. Both of the previous approaches aim to construct foundation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

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

3. Learning 4D Embodied World Models

We introduce the 4D Embodied World Model, which predicts future RGB, depth, and normal maps based on a given input image and text. We leverage a pre-trained video diffusion model as the backbone to predict this rich set of 2D geometric information. Then, we propose an efficient method to

113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140

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

147 convert the RGB-DN video into 4D scenes.

148

3.1. 4D Embodied Video Data Annotation

149 Learning 4D embodied world models requires large-scale
150 4D datasets, which are expensive to collect in the real world.
151 In this section, we present a data annotation pipeline that
152 enables us to automatically construct 4D datasets from existing
153 video datasets. For an illustration of this process, see Fig. 2.154 Given an input video \mathcal{V} , our data annotation protocol
155 aims to automatically obtain 4D annotations of depth map se-
156 quence \mathcal{D} and normal map sequence \mathcal{N} . Video depth/normal
157 requires consistent predictions between consecutive frames
158 [4, 43]. Therefore, we utilize flow as guidance [39]. By
159 predicting the optical flow, we can guarantee that static
160 areas of a scene remain temporally consistent between
161 consecutive frames. Note that the optical flow method
162 is effective in this case, as existing manipulation datasets
163 [5, 6, 18, 27, 35, 36, 46, 57] are based on fixed camera poses.164 To be specific, we employ RAFT [55], an efficient opti-
165 cal flow estimation model, to generate the optical flow \mathcal{F} .
166 The optical flow between consecutive frames is computed
167 as $\mathcal{F} = \text{RAFT}(\mathcal{V})$. The next step involves generating the
168 3D annotations: depth \mathcal{D} and surface normal \mathcal{N} . These 3D
169 annotations provide richer visual context and enhance the
170 understanding of the 4D environment represented by \mathcal{V} . Re-
171 cent advances in monocular depth and normal estimators,
172 typically based on models like Vision Transformers (ViT) or
173 Diffusion Models, have shown strong generalization capabili-
174 ties across diverse datasets.175 Then, we leverage the diffusion model Marigold [30]
176 to estimate \mathcal{D}^i and \mathcal{N}^i from each video frame \mathcal{V}^i . The
177 frame is first encoded into a latent representation \mathbf{v}^i using a
178 variational autoencoder (VAE) [33, 51]: $\mathbf{v}^i = \text{Encoder}(\mathcal{V}^i)$.
179 At the same time, an initial depth or normal latent map \mathbf{z}_T^i is
180 sampled from a Gaussian distribution: $\mathbf{z}_T^i \sim \mathcal{N}(0, 1)$. Then
181 we iteratively apply the learned denoiser U-Net [53] ϵ on
182 each timestep t to reconstruct the \mathbf{z}_0^i and \mathcal{Z}^i :

183
$$\mathbf{z}_{t-1}^i = \epsilon(\mathbf{v}^i, \mathbf{z}_t^i, t) \quad \text{and} \quad \mathcal{Z}^i = \text{Decoder}(\mathbf{z}_0^i) \quad (1)$$

184 where $(\mathbf{z}, \mathcal{Z}) \in \{(\mathbf{d}, \mathcal{D}), (\mathbf{n}, \mathcal{N})\}$.185 However, most existing approaches [3, 30, 63] are pri-
186 marily designed for image-based inputs and often struggle187 when processing videos, resulting in unstable predictions. To
188 overcome this limitation, we propose a novel technique that
189 leverages optical flow as guidance to refine depth estimation.
190 The key insight is that optical flow can capture how static
191 backgrounds and objects move within a scene. Therefore,
192 during the inference stage of the diffusion models, we intro-
193 duce a loss function that enforces consistency between the
194 backgrounds of consecutive frames. Specifically, we define
195 a static region mask for the i -th frame based on the optical
196 flow magnitude: $\mathcal{M}^i = (\|\mathcal{F}^i\| < c)$ where c is a predefined
197 threshold. In practice, we erode the mask to ensure it covers
198 a robust region. Finally, we define the loss function and
199 integrate its gradient into Eq. 1 for every denoising step:

200
$$\mathcal{L}(\mathbf{z}_t^i) = \left\| \text{Decoder}\left(\epsilon(\mathbf{v}^i, \mathbf{z}_t^i, t)\right) \circ (\mathcal{M}^i \cap \mathcal{M}^{i-1}) \right. \\ \left. - \mathcal{Z}^{i-1} \circ (\mathcal{M}^i \cap \mathcal{M}^{i-1}) \right\|^2 \quad (2)$$

201
$$\mathbf{z}_{t-1}^i = \epsilon_\theta \left[\mathbf{v}^i, \left(\mathbf{z}_t^i - w \nabla_{\mathbf{z}_t^i} \mathcal{L} \right), t \right] \quad (3)$$

202 where w is a guidance weight, \circ represents the element-
203 wise product and $\mathcal{M}_i \cap \mathcal{M}_{i-1}$ selects the overlapping stable
204 regions between two consecutive frames. The pseudocode is
205 presented in the appendix.206

3.2. Preliminaries on Video Diffusion Models

207 Diffusion models [24, 52] are capable of learning the data
208 distribution $p(x)$ by progressively adding noise to the data
209 until it resembles a Gaussian distribution through a forward
210 process. During inference, a denoiser ϵ is trained to recover
211 the data from this noisy state. Latent video diffusion models
212 [70] utilize a Variational Autoencoder (VAE) [33, 56], in the
213 latent space of the data, maintaining high-quality outputs
214 while more efficiently modeling the data distribution. In this
215 section, we formulate the task of RGB \mathcal{V} , depth \mathcal{D} , and nor-
216 mal \mathcal{N} map video generation as a conditional denoising gen-
217 eration task, i.e., we model the distribution $p(\mathbf{v}, \mathbf{d}, \mathbf{n} | \mathbf{v}_0, \mathcal{T})$,
218 where $\mathbf{v}, \mathbf{d}, \mathbf{n}$ represent the predicted future latent sequences
219 of RGB, depth, and normal maps, respectively. The con-
220 dition \mathbf{v}_0 is a given RGB image latent, and \mathcal{T} denotes the
221 instruction provided by the user.

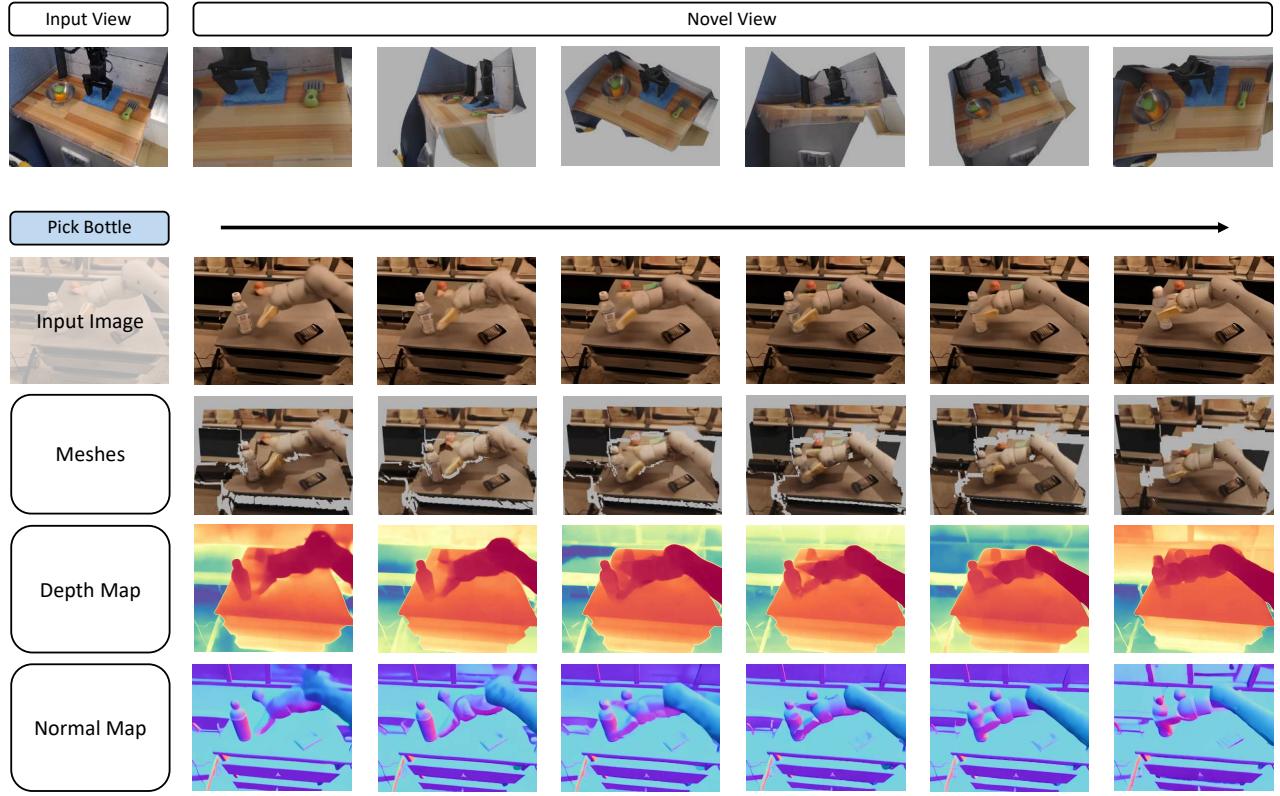


Figure 3. 4D Reconstruction and Video Generation. The first row shows the results of our 4D reconstruction method from multiple views on a single frame. The following rows present the inference results of our method given the input image and instruction, along with the reconstructed meshes.

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

$$225 \quad q(\mathbf{z}_t | \mathbf{z}_{t-1}) = \mathcal{N}(\mathbf{z}_t; \sqrt{\alpha_t} \mathbf{z}_{t-1}, (1 - \alpha_t) \mathbf{I}) \quad (4)$$

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

$$234 \quad p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{v}_0, \mathcal{T}) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t, \mathbf{v}_0, \mathcal{T}), \Sigma_\theta(\mathbf{x}_t, t)) \quad (5)$$

235 where

$$236 \quad \mu_\theta(\mathbf{x}_t, t, \mathbf{v}_0, \mathcal{T}) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \epsilon_\theta(\mathbf{x}_t, t, \mathbf{v}_0, \mathcal{T}) \right). \quad (6)$$

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

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

$$246 \quad L = \mathbb{E}_{\mathbf{v}_0, \mathcal{T}, t, \epsilon} \left[\left\| [\epsilon_v, \epsilon_d, \epsilon_n] - \epsilon_\theta(\mathbf{x}_t, t, \mathbf{v}_0, \mathcal{T}) \right\|^2 \right] \quad (7)$$

3.3. RGB, Depth and Normal Video Predictions

247 Training a diffusion model to model temporal RGB-DN data
248 is a challenging task. To effectively train RGB video models,
249 large-scale video datasets with billions of high-quality sam-
250 ples are used [70]. In contrast, even through automatic an-
251 notation, our dataset of RGB-DN data contains only around
252 one million data points, which is insufficient to train a world
253 model from scratch. To address this, we finetune the Open-
254 Sora [70] as our RGB-DN prediction model and directly
255 leverage the pre-trained knowledge inside the model to ef-
256 fectively bootstrap our 4D model.

257 To implement this, we use the temporal VAE [33, 56]
258 in Open-Sora [70] to separately encode RGB, depth, and
259 normal images for each frame of a video, without additional
260 fine-tuning of the VAE. We then expand Open-Sora STDiT's
261

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

269 3.4. 4D Scene Reconstruction from RGB-DN Video

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

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

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

$$311 \min_{\tilde{d}} \iint_{\Omega} (\tilde{n}_z \partial_u \tilde{d} + n_x)^2 + (\tilde{n}_z \partial_v \tilde{d} + n_y)^2 dudv. \quad (8)$$

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

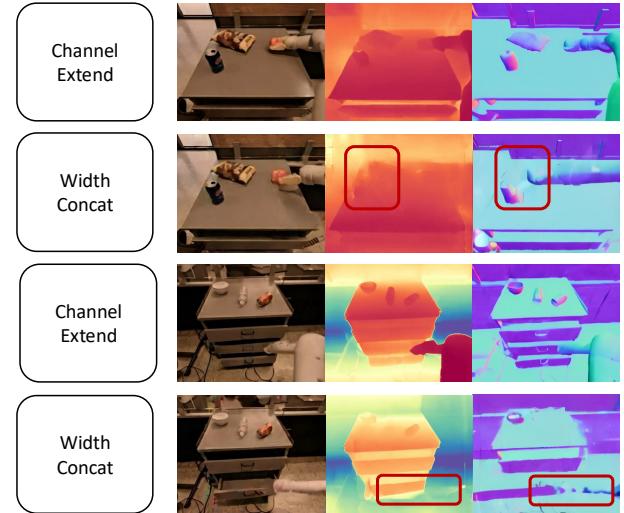


Figure 4. Channel concatenation improves visual quality and ensures consistency across maps.

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

$$316 \tilde{d}_{t+1} = \arg \min_{\tilde{d}} (A\tilde{d} - b)^T W(\tilde{d}_t) (A\tilde{d} - b) \stackrel{\text{def}}{=} \arg \min_{\tilde{d}} \mathcal{L}(\tilde{\mathcal{D}}, \mathcal{N}^i), \quad (9)$$

317 where A and b are defined by predicted normals and camera
 318 intrinsics.

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

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

$$341 \mathcal{L}_d(\tilde{\mathcal{D}}, \hat{\mathcal{D}}^{i-1}, M_d^i, \mathcal{F}^i, \mathcal{F}^{i-1}) = \left\| \tilde{\mathcal{D}}^i \circ M_d^i - \mathcal{D}^{i \rightarrow (i-1)} \circ M_d^i \right\|^2. \quad (10)$$

342 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	0.0945	0.2022

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

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

346
$$\mathcal{L}_g(\mathcal{D}_1, \mathcal{D}_2, \mathcal{M}) = \|\mathcal{D}_1 \circ \mathcal{M} - \mathcal{D}_2 \circ \mathcal{M}\|^2 \quad (11)$$

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

355 The overall loss objective we optimize is given by:

356
$$\arg \min_{\tilde{\mathcal{D}}} \mathcal{L}(\tilde{\mathcal{D}}, \mathcal{N}^i) + \lambda_b \mathcal{L}_g(\tilde{\mathcal{D}}, \hat{\mathcal{D}}^{i-1}, \mathcal{M}_b^i) + \quad (12)$$

357
$$+ \lambda_d \mathcal{L}_d(\tilde{\mathcal{D}}, \hat{\mathcal{D}}^{i-1}, \mathcal{M}_d^i, \mathcal{F}^i, \mathcal{F}^{i-1}) + \quad (13)$$

358
$$+ \lambda_{g2} \mathcal{L}_g(\tilde{\mathcal{D}}, \mathcal{D}^i, \mathcal{M}_b^i) + \lambda_{g1} \mathcal{L}_g(\tilde{\mathcal{D}}, \mathcal{D}^i, \mathcal{M}_d^i). \quad (14)$$

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

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

375 3.5. Inverse Dynamics Models from 4D Scenes

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

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

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.

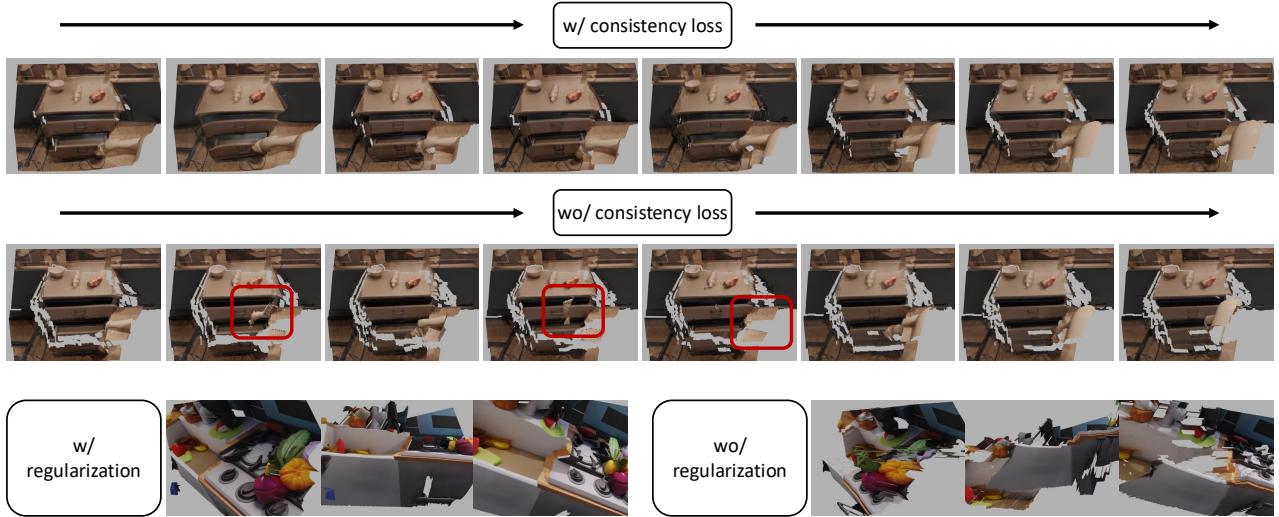


Figure 5. Effect of Consistency and Regularization Losses on 4D Mesh Reconstruction. The red boxes highlight the inconsistent regions.

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	0.7382	66.67	3.61	~2 hours
Ours	12.99	0.4262	0.6051	83.02	3.73	~ 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 ↑		δ_1 ↑	δ_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	✗	25.78	71.89	21.86	16.14	76.59	91.54	26.59	15.73	24.36
16	✓	25.45	74.98	21.94	16.53	77.13	92.15	16.23	7.78	38.50
										32.98

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	57.1
Ours	95.2	28.6	23.8	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.

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, Xiaojian 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

- 673 [29] Yunfan Jiang, Agrim Gupta, Zichen Zhang, Guanzhi Wang,
674 Yongqiang Dou, Yanjun Chen, Li Fei-Fei, Anima Anand-
675 kumar, Yuke Zhu, and Linxi Fan. Vima: General robot
676 manipulation with multimodal prompts. *arXiv preprint*
677 *arXiv:2210.03094*, 2(3):6, 2022. 2
- 678 [30] Bingxin Ke, Anton Obukhov, Shengyu Huang, Nando Met-
679 zger, Rodrigo Caye Daudt, and Konrad Schindler. Repurposing
680 diffusion-based image generators for monocular depth
681 estimation. In *Proceedings of the IEEE/CVF Conference on*
682 *Computer Vision and Pattern Recognition*, pages 9492–9502,
683 2024. 2, 3, 5, 6, 8
- 684 [31] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and
685 George Drettakis. 3d gaussian splatting for real-time radiance
686 field rendering. *ACM Trans. Graph.*, 42(4):139–1, 2023. 7
- 687 [32] Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted
688 Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailev, Ethan
689 Foster, Grace Lam, Pannag Sanketi, et al. Openvla: An
690 open-source vision-language-action model. *arXiv preprint*
691 *arXiv:2406.09246*, 2024. 2
- 692 [33] Diederik P Kingma. Auto-encoding variational bayes. *arXiv*
693 *preprint arXiv:1312.6114*, 2013. 3, 4
- 694 [34] Diederik P Kingma. Auto-encoding variational bayes. *arXiv*
695 *preprint arXiv:1312.6114*, 2013. 8
- 696 [35] Danica Kragic, Mårten Björkman, Henrik I Christensen, and
697 Jan-Olof Eklundh. Vision for robotic object manipulation in
698 domestic settings. *Robotics and autonomous Systems*, 52(1):
699 85–100, 2005. 3
- 700 [36] Nikolaos Kyriazis and Antonis Argyros. Physically plausible
701 3d scene tracking: The single actor hypothesis. In *Proceed-
702 ings of the IEEE conference on computer vision and pattern
703 recognition*, pages 9–16, 2013. 3
- 704 [37] LAION-AI. Aesthetic predictor. [https://github.com/
705 LAION-AI/aesthetic-predictor](https://github.com/LAION-AI/aesthetic-predictor), 2022. 7
- 706 [38] Timothée Lesort, Natalia Díaz-Rodríguez, Jean-François
707 Goudou, and David Filliat. State representation learning
708 for control: An overview. *Neural Networks*, 108:379–392,
709 2018. 2
- 710 [39] Siyuan Li, Yue Luo, Ye Zhu, Xun Zhao, Yu Li, and Ying Shan.
711 Enforcing temporal consistency in video depth estimation. In
712 *Proceedings of the IEEE/CVF International Conference on*
713 *Computer Vision*, pages 1145–1154, 2021. 3
- 714 [40] Shuang Li, Xavier Puig, Chris Paxton, Yilun Du, Clinton
715 Wang, Linxi Fan, Tao Chen, De-An Huang, Ekin Akyürek,
716 Anima Anandkumar, et al. Pre-trained language models for
717 interactive decision-making. *Advances in Neural Information
718 Processing Systems*, 35:31199–31212, 2022. 2
- 719 [41] Junbang Liang, Ruoshi Liu, Ege Ozguroglu, Sruthi Sudhakar,
720 Achal Dave, Pavel Tokmakov, Shuran Song, and Carl Von-
721 drick. Dreamitate: Real-world visuomotor policy learning via
722 video generation. *arXiv preprint arXiv:2406.16862*, 2024. 1,
723 2
- 724 [42] Lennart Ljung and Torkel Glad. *Modeling of dynamic systems*.
725 Prentice-Hall, Inc., 1994. 2
- 726 [43] Fei Luo, Lin Wei, and Chunxia Xiao. Stable depth estimation
727 within consecutive video frames. In *Advances in Computer
728 Graphics: 38th Computer Graphics International Conference,
729 CGI 2021, Virtual Event, September 6–10, 2021, Proceedings*
730 38, pages 54–66. Springer, 2021. 3
- 731 [44] Vincent Micheli, Eloi Alonso, and François Fleuret. Trans-
732 formers are sample efficient world models. *arXiv preprint*
733 *arXiv:2209.00588*, 2022. 2
- 734 [45] Alex Nichol, Heewoo Jun, Prafulla Dhariwal, Pamela
735 Mishkin, and Mark Chen. Point-e: A system for generat-
736 ing 3d point clouds from complex prompts. *arXiv preprint*
737 *arXiv:2212.08751*, 2022. 6
- 738 [46] Antje Nuthmann and Teresa Canas-Bajo. Visual search in
739 naturalistic scenes from foveal to peripheral vision: A compar-
740 ision between dynamic and static displays. *Journal of Vision*,
741 22(1):10–10, 2022. 3
- 742 [47] OpenAI. Video generation models as world simu-
743 lators. [https://openai.com/index/video-
744 generation-models-as-world-simulators/](https://openai.com/index/video-generation-models-as-world-simulators/), 2023. Accessed: 2024-10-01. 2
- 745 [48] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J
746 Guibas. Pointnet++: Deep hierarchical feature learning on
747 point sets in a metric space. *Advances in neural information
748 processing systems*, 30, 2017. 6
- 749 [49] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya
750 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,
751 Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning
752 transferable visual models from natural language supervi-
753 sion. In *International conference on machine learning*, pages
754 8748–8763. PMLR, 2021. 6
- 755 [50] Shreyas Sundara Raman, Vanya Cohen, Eric Rosen, Ifrah
756 Idrees, David Paulius, and Stefanie Tellex. Planning with
757 large language models via corrective re-prompting. *arXiv*
758 *preprint arXiv:2211.09935*, 2022. 2
- 759 [51] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wier-
760 stra. Stochastic backpropagation and approximate inference
761 in deep generative models. In *International conference on*
762 *machine learning*, pages 1278–1286. PMLR, 2014. 3
- 763 [52] Robin Rombach, Andreas Blattmann, Dominik Lorenz,
764 Patrick Esser, and Björn Ommer. High-resolution image
765 synthesis with latent diffusion models. In *Proceedings of*
766 *the IEEE/CVF conference on computer vision and pattern
767 recognition*, pages 10684–10695, 2022. 3
- 768 [53] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-
769 net: Convolutional networks for biomedical image segmen-
770 tation. In *Medical image computing and computer-assisted
771 intervention—MICCAI 2015: 18th international conference,
772 Munich, Germany, October 5–9, 2015, proceedings, part III
773 I8*, pages 234–241. Springer, 2015. 3
- 774 [54] Richard S Sutton. Dyna, an integrated architecture for learn-
775 ing, planning, and reacting. *ACM Sigart Bulletin*, 2(4):160–
776 163, 1991. 2
- 777 [55] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field
778 transforms for optical flow. In *Computer Vision—ECCV 2020:
779 16th European Conference, Glasgow, UK, August 23–28,
780 2020, Proceedings, Part II 16*, pages 402–419. Springer, 2020.
781 3, 5
- 782 [56] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete
783 representation learning. *Advances in neural information pro-
784 cessing systems*, 30, 2017. 3, 4, 8
- 785 [57] Homer Walke, Kevin Black, Abraham Lee, Moo Jin Kim,
786 Max Du, Chongyi Zheng, Tony Zhao, Philippe Hansen-
787 Estruch, Quan Vuong, Andre He, Vivek Myers, Kuan Fang,
788 10

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