

Learning 4D Embodied World Models

Supplementary Material

001 A. Implementation Details

002 A.1. Data Annotation Details

003 For data generation, we use the Marigold [4] DDIM Scheduler,
 004 performing 100 denoising steps in total. During the last 30 steps, optical flow guide inference was incorporated.
 005 The threshold for background mask computation was set to 0.5, and the gradient descent weight was configured as 2000.
 006 The pseudocode for the guidance stage is as follows:

Algorithm 1 Denoising With Optical Flow Consistency

```

1: Parameters:  $w$            ▷ Gradient descent weight
2: Inputs:
   •  $\mathbf{z}_t^i$ : Noisy depth/normal latent variable at timestep  $t$ 
   •  $\mathbf{v}^i$ : Latent encoding of the video frame  $V^i$ 
   •  $\mathcal{M}^i, \mathcal{M}^{i-1}$ : Static region masks for the current and previous frames
   •  $\mathcal{Z}^{i-1}$ : Depth/normal annotation from the previous frame
   •  $t$ : Current timestep in the denoising process
3: Outputs:  $\mathbf{z}_{t-1}^i$ : Noisy images at timestep  $t - 1$ 
4: function DENOISINGONESTEP( $\mathbf{z}^i_t$ )
5:    $\mathbf{z}_t^i$ : require grads  $\leftarrow$  True
6:    $\mathbf{z} = \epsilon_\theta(\mathbf{v}^i, \mathbf{z}_t^i, t)$            ▷ First denoise
7:    $\mathcal{L} = \|\text{Decoder}(z) \circ (\mathcal{M}^i \cap \mathcal{M}^{i-1}) - \mathcal{Z}^{i-1} \circ (\mathcal{M}^i \cap \mathcal{M}^{i-1})\|^2$       ▷ By Eq.2
8:    $\mathcal{L}.\text{backward}()$           ▷ To get the value of  $\nabla_{\mathbf{z}_t^i} \mathcal{L}$ 
9:    $\mathbf{z}_t^i \leftarrow \mathbf{z}_t^i - w \nabla_{\mathbf{z}_t^i} \mathcal{L}$ 
10:   $\mathbf{z}_{t-1}^i = \epsilon_\theta(\mathbf{v}^i, \mathbf{z}_t^i, t)$            ▷ Second denoise
11: end function

```

009 Figure 1 illustrates the denoising process within the diffusion model, which is repeated at each timestep. Each step
 010 consists of three stages: **First Denoise**, where the noisy lat-
 011 ent variable is passed through the Marigold for an initial
 012 forward; **Loss Backward**, where the loss is computed using
 013 static region masks and depth from the previous frame,
 014 followed by gradient-based updates to refine the latent variable;
 015 and **Second Denoise**, where the adjusted latent variable is
 016 further processed to produce a refined output. This iterative
 017 process ensures consistency across timesteps while improv-
 018 ing depth/normal estimation.

020 A.2. Video Diffusion Model Details

021 Our RGB-DN video diffusion model's detailed training proto-
 022 colog is shown in Figure 2. For the RGB-DN video latent, we
 023 add noise to all frames except the first one, condition it on
 024 text instructions, and predict the added noise. We trained our

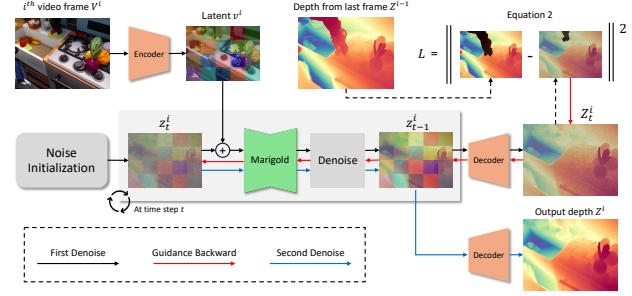


Figure 1. **Data generation pipeline**, comprising first denoise, guidance backward, and second denoise stages to produce refined depth/normal estimations.

video diffusion model using the STDiT3-XL/2 architecture, fine-tuned on OpenSora-v1.2 [10], employing 6×8 V100 GPUs. The model processes videos with 16 frames, using gradient checkpointing to optimize memory usage. We set a batch size of 2 for acceleration, used bf16 precision, and applied the ZeRO2 [7] optimization plugin. Additionally, we leveraged a T5 text encoder [6] for conditioning. For sampling, we use the RFlow scheduler with logit-normal sampling across 30 steps and set a classifier-free guidance scale of 7.0.

The training spanned 40,000 iterations with an initial learning rate 1e-4, 1.0 gradient clipping, and a 1,000-step warmup. The optimizer incorporated Adam with epsilon set to 1e-15, and an exponential moving average (EMA [5]) decay of 0.99 was used to stabilize training.

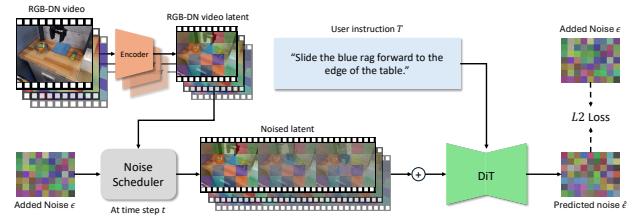


Figure 2. Training Pipeline of the RGB-DN Video Diffusion Model

A.3. 4D Meshes Generation

The parameters for the loss term in Eq.12 are set differently for the RT-1 [2] and RLBench [3] datasets, as shown in the table below:

In Figure 4, we present a visualization of the 3D robotic scene reconstruction optimized using our proposed method in the BridgeV2 [8] dataset. After estimating the depth and normal with the estimator, we refine the outputs to reconstruct the scene accurately. The figure includes untextured

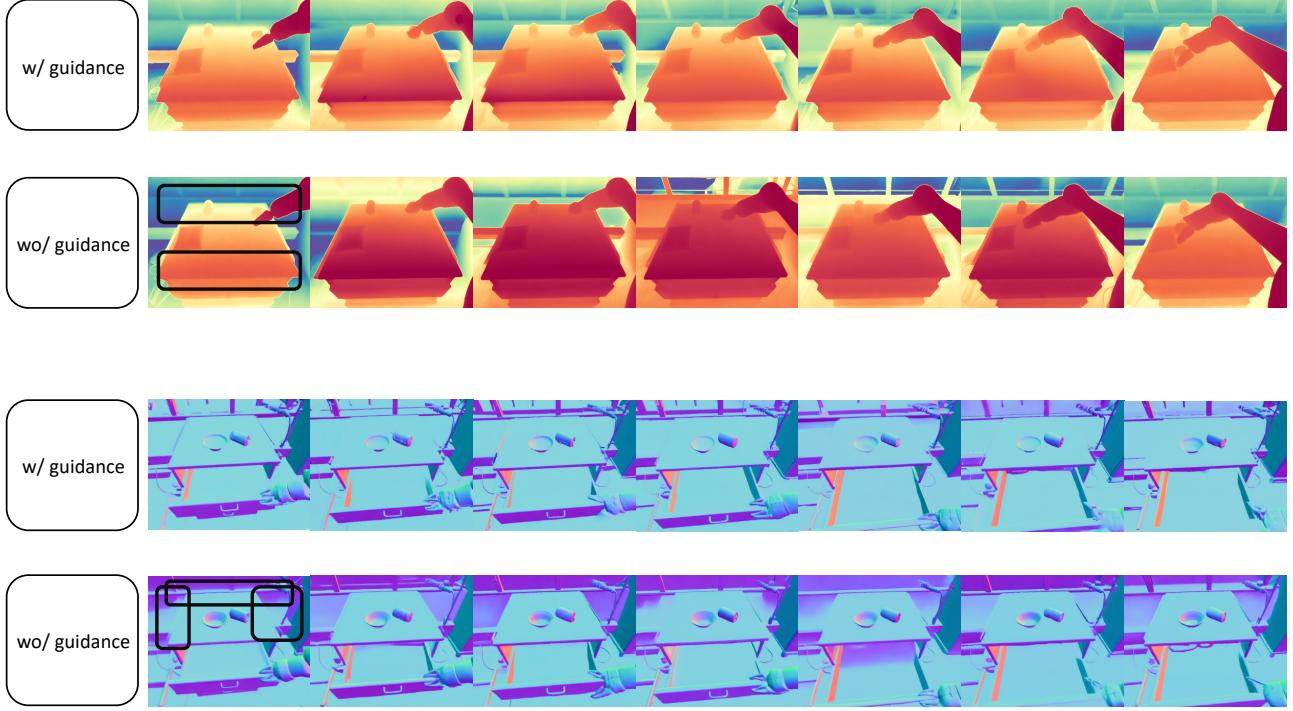


Figure 3. Qualitative Results of Data Generation Methods with and without Guidance, where the black boxes highlight areas of inconsistency.

Dataset	λ_d	λ_b	λ_{g1}	λ_{g2}
RT-1	20	200	20	20
RLBench	20	200	2	2

Table 1. Loss Term Parameters for RT-1 and RLBench Datasets

049 rendering and texture-rendered views, where the wall tex-
 050 tures are significantly enhanced due to normal optimization.
 051 The side perspective view shows the improved shape and
 052 geometry reconstruction. Notably, the wall and table surfaces
 053 are well-aligned, appearing perpendicular to each other, fur-
 054 ther validating the effectiveness of our optimization process
 055 in capturing accurate spatial relationships.

A.4. Implementation Details for Robotics Planning

057 For the RLBench training, we adopted the same architecture
 058 and methods as our video diffusion model, with the primary
 059 difference being that we used 13 frames and fine-tuned the
 060 model on the RT-1 dataset.

061 For the action prediction stage, we first filter out the back-
 062 ground and floor from the data, focusing only on the points
 063 of the table and the objects manipulated by the robotic arm,
 064 and then sample 1024 points from the filtered point cloud.
 065 In our inverse dynamic model, the PointNet extracts features
 066 from this point cloud, concatenated with the instruction's

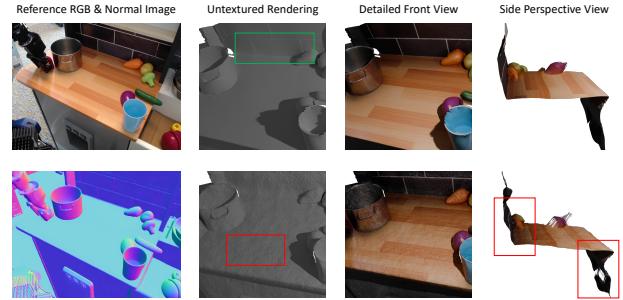


Figure 4. Visualization of the optimized 3D robotic scene reconstruction using our method. The untextured renderings show enhanced detail (green box) and improved surface smoothness (red box). The side perspective view highlights accurate shape and geometry optimization, including the perpendicular alignment of the wall and table (red boxes).

language embedding, and passed into a 4-layer MLP, finally
 067 outputting the 7DoF actions.
 068

B. More Experiments

B.1. Data Annotation

071 In this section, we first compare our data generation method
 072 with 3D-VLA [9]. They use ZoeDepth [1] for depth map
 073 estimation and directly map them into 3D space. The com-

Pretrained	Channel	FVD \uparrow	RGB SSIM \uparrow	PSNR \uparrow	Depth AbsRel \downarrow	$\delta_1 \uparrow$	$\delta_2 \uparrow$	Normal Mean \downarrow	Median \downarrow	11.25° \uparrow	Consistency Edge-sim \uparrow
\times	\times	60.32	56.44	16.02	37.65	30.33	71.35	128.57	139.50	0.0	50.52
\checkmark	\times	55.04	56.32	16.36	34.39	34.59	67.34	128.77	139.61	0.0	51.62
\times	\checkmark	54.11	65.90	19.28	18.40	65.02	91.20	12.94	7.58	25.02	53.39
\checkmark	\checkmark	56.82	66.32	19.64	15.35	77.54	93.62	11.06	5.93	43.59	57.26

Table 2. Ablation study of video prediction task on RLBench simulation. “Pretrained” \checkmark indicates fine-tuning models pre-trained on the RT-1 dataset. “Channel” \checkmark indicates concatenating RGB, depth, and normal maps along the channel dimension, while \times refers to concatenation along the width.

parison results, shown in Figure 5, evaluate the quality of point cloud generation for both methods, with cubes replacing vertices for rendering. Our generated data demonstrates higher realism, while 3D-VLA exhibits noticeable shape distortion.

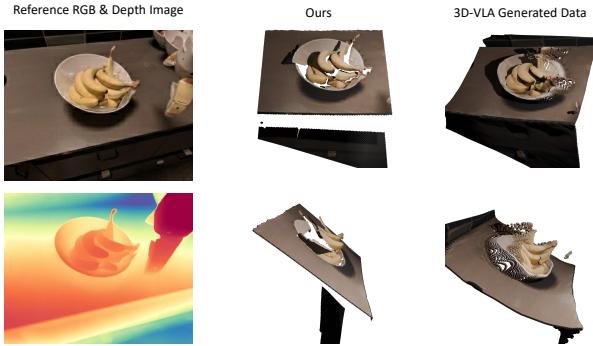


Figure 5. Comparison of point cloud generation quality between our method and 3D-VLA

We also show an ablation study comparing data generation methods with and without guidance. As shown in Figure 3, we focus on the qualitative results to evaluate the impact of guidance. The black boxes in the figure highlight areas of inconsistency, demonstrating that the Guided Marigold method achieves superior frame-to-frame coherence compared to the unguided Marigold-LCM. This emphasizes the importance of guidance in producing consistent dynamic scenes.

B.2. RGB-DN Prediction

We conducted an ablation study of the model framework using the RLBench simulation environment, as it provides ground truth depth values, making the evaluation more convincing. Since RLBench does not natively offer normal map outputs, we employed guided Marigold to generate normal annotations. The metrics used in this table are consistent with those described in the main text, ensuring a direct comparison with the results reported therein.

As shown in Table 2, concatenating RGB, depth, and normal maps along the width dimension performs poorly,

particularly on normal map-related metrics, indicating that this approach is ineffective for normal map processing. In contrast, concatenating these inputs along the channel dimension significantly improves performance across all tasks. Furthermore, fine-tuning pre-trained models on the RT-1 dataset further enhances results. Combining pre-trained models and channel-wise concatenation yields the best overall performance.

B.3. More Qualitative Results

Figure 7 showcases qualitative results from the RT-1 dataset, illustrating our approach’s capabilities in generating multi-modal visual representations. Our video diffusion model directly produces the RGB, depth, and normal maps, while the meshes are rendered from the reconstructed outputs. These results highlight the robustness of our method across different visual modalities. *Additional video results* can be found in the supplementary materials folder for further analysis and evaluation.

B.4. Action Planning

One potential application of our generated mesh is to extract action trajectories directly. As illustrated in Figure 6, we track the robotic arm in the video to capture its motion path. This trajectory is subsequently lifted into 3D space, enabling the reconstruction of the robot arm’s action trajectory. The red line in the visualization represents the captured action trajectory.



Figure 6. Tracking and visualization of robotic arm action trajectories on the Bridge dataset

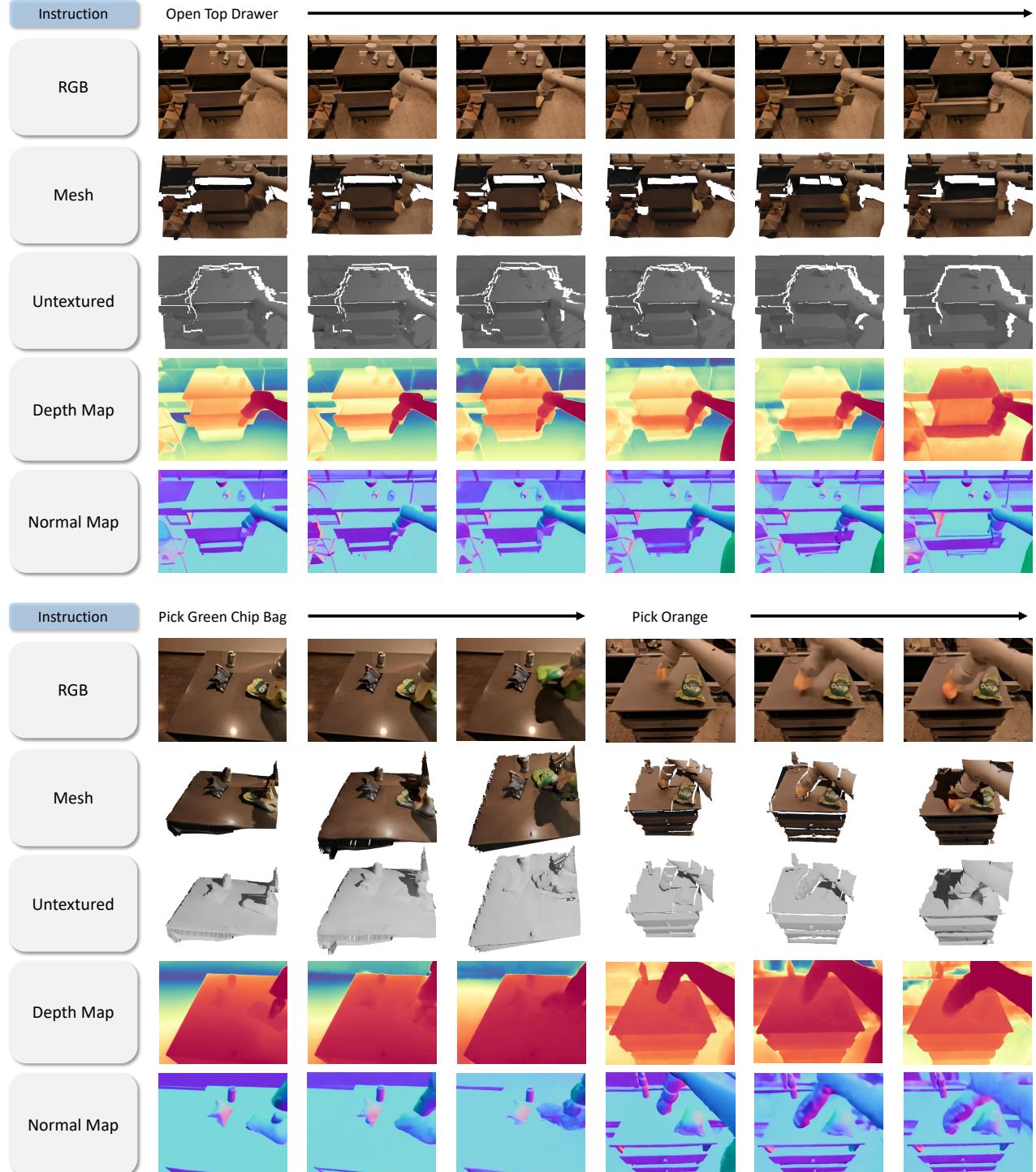


Figure 7. More qualitative results

125

References

- 126 [1] Shariq Farooq Bhat, Reiner Birkl, Diana Wofk, Peter
127 Wonka, and Matthias Müller. Zoedepth: Zero-shot trans-
128 fer by combining relative and metric depth. *arXiv preprint*
129 *arXiv:2302.12288*, 2023. 2
- 130 [2] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen
131 Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakr-
132 ishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al.
133 Rt-1: Robotics transformer for real-world control at scale.
134 *arXiv preprint arXiv:2212.06817*, 2022. 1
- 135 [3] Stephen James, Zicong Ma, David Rovick Arrojo, and An-
136 drew J. Davison. Rlbench: The robot learning benchmark &
137 learning environment. *IEEE Robotics and Automation Letters*,
138 2020. 1
- 139 [4] Bingxin Ke, Anton Obukhov, Shengyu Huang, Nando Met-
140 zger, Rodrigo Caye Daudt, and Konrad Schindler. Repurpos-
141 ing diffusion-based image generators for monocular depth
142 estimation. In *Proceedings of the IEEE/CVF Conference on*
143 *Computer Vision and Pattern Recognition*, pages 9492–9502,
144 2024. 1
- 145 [5] Frank Klinker. Exponential moving average versus moving
146 exponential average. *Mathematische Semesterberichte*, 58:
147 97–107, 2011. 1
- 148 [6] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee,
149 Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and
150 Peter J Liu. Exploring the limits of transfer learning with a
151 unified text-to-text transformer. *Journal of machine learning*
152 *research*, 21(140):1–67, 2020. 1
- 153 [7] Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yux-
154 iong He. Zero: Memory optimizations toward training trillion
155 parameter models. In *SC20: International Conference for*
156 *High Performance Computing, Networking, Storage and Anal-*
157 *ysis*, pages 1–16. IEEE, 2020. 1
- 158 [8] Homer Walke, Kevin Black, Abraham Lee, Moo Jin Kim,
159 Max Du, Chongyi Zheng, Tony Zhao, Philippe Hansen-
160 Estruch, Quan Vuong, Andre He, Vivek Myers, Kuan Fang,
161 Chelsea Finn, and Sergey Levine. Bridgedata v2: A dataset
162 for robot learning at scale. In *Conference on Robot Learning*
163 (*CoRL*), 2023. 1
- 164 [9] Haoyu Zhen, Xiaowen Qiu, Peihao Chen, Jincheng Yang, Xin
165 Yan, Yilun Du, Yining Hong, and Chuang Gan. 3d-vla: A
166 3d vision-language-action generative world model. *arXiv*
167 *preprint arXiv:2403.09631*, 2024. 2
- 168 [10] Zangwei Zheng, Xiangyu Peng, Tianji Yang, Chenhui Shen,
169 Shenggui Li, Hongxin Liu, Yukun Zhou, Tianyi Li, and Yang
170 You. Open-sora: Democratizing efficient video production
171 for all, 2024. 1