
Vision-Language Models Meet Video Diffusion: Plug-and-Play World Models for Scalable Embodied Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 World models have emerged as powerful tools for planning and decision-making
2 in embodied agents, especially with recent advances in diffusion-based video gen-
3 eration that enable high-fidelity future prediction. However, challenges in scaling
4 high-resolution, instruction-following video generation and bridging simulated
5 predictions to real-world robot control hinder practical deployment. To address
6 these, we introduce a novel framework that repurposes video diffusion models
7 as zero-shot policy generators conditioned on abstract action primitives dynam-
8 ically derived from language instructions via large vision-language models. By
9 extracting 6-DoF end-effector trajectories from predicted video frames and en-
10 abling closed-loop execution with observation feedback, our method maintains
11 interpretability while avoiding error accumulation. This modular approach lever-
12 ages the inductive biases of video generation and vision-language reasoning, al-
13 lowing plug-and-play generalization across tasks. Furthermore, it enables policy
14 distillation from unannotated visual demonstrations—ranging from simulated roll-
15 outs to real-world or web-sourced videos—helping mitigate data scarcity in embed-
16 ded learning.

17

1 Introduction

18 Embodied agents capable of planning and decision-making in complex environments rely heavily
19 on internal representations of the world, commonly referred to as *world models*. Recent advances in
20 video generation—particularly diffusion-based models—have enabled high-fidelity visual future pre-
21 diction, revitalizing interest in world model-centric architectures. However, two critical challenges
22 persist: 1) Scalability-Complexity Trade-off: High-resolution, long-horizon video generation with
23 fine-grained instruction following demands immense computational overheads, limiting real-time
24 deployment; 2) Zero-Shot Policy Gap: Bridging simulated visual futures to real-world robot actions
25 without task-specific adaptation remains an open problem. On the other hand, while hierarchical ar-
26 chitectures have demonstrated success in end-to-end vision-language-action models by decoupling
27 high-level planning from low-level action prediction (e.g., SayCan), such structured designs remain
28 relatively underexplored in the context of world models. Furthermore, end-to-end paradigms—such
29 as RT-1/RT-2—that directly map vision and language to actions entangle high-level semantic reason-
30 ing with low-level motor/end-effector dynamics. This design burdens a single model with multiple
31 abstraction levels, leading to poor data efficiency, limited explainability, and weak generalization
32 across tasks or morphologies.

33 To address these limitations, we propose a modular and hierarchical framework that reinterprets
34 video diffusion models as **primitive-level world models** for zero-shot policy generation. Rather
35 than attempting to learn task-level action distributions directly, our method decomposes complex

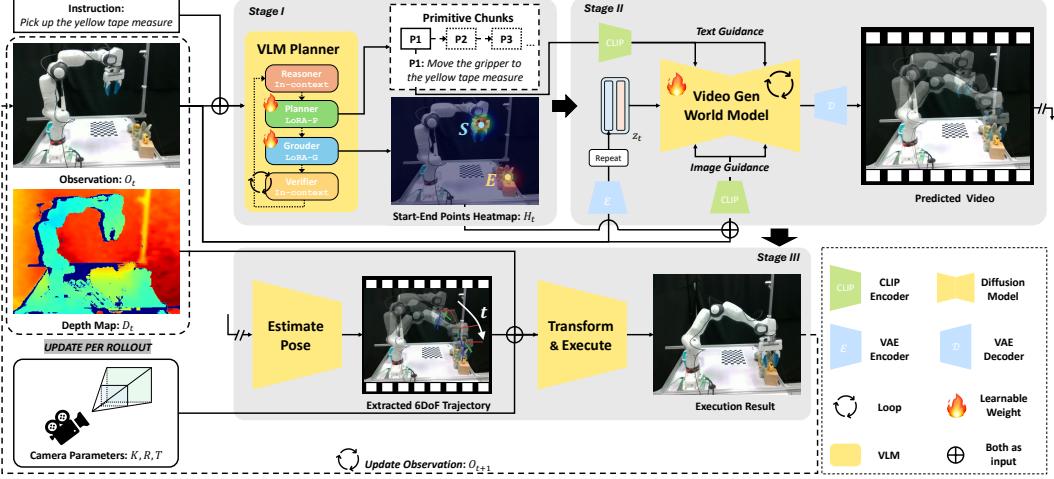


Figure 1: This workflow diagram illustrates a video generation framework structured into three stages: Stage 1 employs a VLM planner, comprising a VLM and two learnable LoRAs, along with a Reasoner and Verifier to ensure the efficiency and accuracy of the planning process. Stage 2 utilizes CLIP and a VAE encoder to generate predicted videos through the Video Generation World Model. Stage 3 focuses on pose estimation for the robotic arm, executing coordinate transformations to produce executable commands for real-world implementation.

36 instructions into semantically aligned atomic units—**action primitives**—and constrains video generation
 37 to short-horizon rollouts grounded in each primitive. For a quick overview of our approach,
 38 please refer to Figure 1.

39 At the heart of our system is a biologically inspired *cerebro-cerebellar hierarchy*: 1) A *cerebral*
 40 planner, instantiated as a vision-language foundation model (VLM), parses natural language and
 41 decomposes high-level tasks into a sequence of interpretable primitives. This narrows the predictive
 42 horizon per rollout, significantly reducing the computational burden otherwise associated with long-
 43 horizon diffusion; 2) A *cerebellar* world model, realized via a lightweight video diffusion model,
 44 predicts visual futures conditioned on each primitive, enabling controllable and interpretable rollout.
 45 This aligns the semantic granularity of visual prediction with the logical structure of instructions,
 46 allowing world models to focus on local dynamics without losing semantic context.

47 To operationalize this hierarchy, we leverage the grounding and spatial reasoning capabilities of the
 48 state-of-the-art VLMs. Given a high-level task instruction and current visual observation, the VLM
 49 first performs semantic decomposition, parsing the instruction into a sequence of executable action
 50 primitives. For the first primitive, the model then identifies two key spatial anchors: the pixel-space
 51 coordinate of the current gripper position (start) and the semantically inferred goal location (end).
 52 These coordinates are rendered as 2D heatmaps, Gaussian-smoothed, and used to produce a spatial
 53 guidance signal. This signal is then provided as a conditioning input to the video diffusion model,
 54 effectively anchoring its generation to a grounded, primitive-specific trajectory plan in pixel space.

55 Critically, our approach ensures that the video diffusion model observes the entire robot arm, in-
 56 cluding the base, joints, and full kinematic chain—unlike prior work that focuses narrowly on the
 57 end-effector or operates in cropped views. This global visibility enables the model to internalize
 58 soft physical constraints, such as joint limits, reachability, and base stability, by learning from visual
 59 data alone. As a result, the generated visual rollouts are not only semantically aligned and visually
 60 plausible, but also physically valid and executable in real environments.

61 Finally, instead of predicting control signals directly, we extract 6-DoF end-effector trajectories
 62 from the video frames using geometry-based visual projection. These trajectories are executed via
 63 Cartesian control, and the resulting observations are fed back to the planner, forming a semi-closed-
 64 loop control pipeline that enhances robustness, avoids error accumulation, and supports plug-and-
 65 play adaptation across embodiments.

66 In summary, our work advances embodied learning along three core axes—**Effectiveness**, **Explainability**, and **Efficiency** (E^3). Our framework effectively combines the semantic parsing strengths of
67 VLMs with the temporal modeling capacity of video diffusion, enabling policy distillation from di-
68 verse visual sources without task-specific retraining. It remains fully explainable through primitive-
69 level decomposition, spatial grounding, and trajectory-based control, offering transparency and
70 safety in real-world deployment. Furthermore, by limiting generation to short-horizon primitives,
71 our system achieves significant efficiency gains over monolithic video models, supporting real-time
72 inference with minimal overhead.
73

74 In particular, our method facilitates improved grounding between natural language and robotic ex-
75 ecution, more efficient data utilization, and soft embodiment-aware constraints through full-arm
76 observation—allowing the model to internalize physically feasible motion dynamics and avoid gener-
77 ating unrealistic or unreachable robot configurations.

78 2 Related Work

79 2.1 Video Generation as World Models

80 Given the overwhelming success of Diffusion Models [78, 36, 80, 79] in image generation [73,
81 72, 25], video diffusion models [50, 37, 35, 52, 8, 77, 87, 16, 17, 90] have emerged as a prominent
82 research focus. With the emergence of Sora [12, 62], Transformer-based video generation has gained
83 traction as a mainstream paradigm. Leveraging the Diffusion Transformer (DiT) architecture [68],
84 these models significantly improve video stability and temporal consistency [101, 55, 97, 46, 71, 98,
85 86]. Empowered by strong visual priors, video diffusion models show promising potential as pixel-
86 level world models. A world model learns to predict future system states, either in a low-dimensional
87 latent space [1, 29, 47] or directly in pixel space [20, 32, 33, 66, 15]. A large body of work [45, 27,
88 95, 12, 6, 89, 13, 38, 102, 23, 54, 21, 38, 51, 81, 69, 100, 103, 34] has framed video generation
89 as a new methodology for learning world models. Latest works, typically DiT-based and following
90 a unified generation paradigm [51, 81, 100], are costly to train and lack real-time responsiveness.
91 While hierarchical approaches [2, 23, 102] offer structure, they struggle with real-time adaptation
92 in dynamic environments and rigid cross-modal alignment. Error accumulation in layered designs
93 remains an open problem. Several approaches [2, 23, 102] adopt hierarchical architectures, among
94 which [2, 23] are the closest to ours. However, these methods lack the integration of vision-language
95 reasoning, planning, and grounding as in our approach. In contrast, our approach is distinguished by
96 its ability to perform zero-shot policy execution, without relying on task-specific trained policies that
97 often suffer from limited generalizability. Furthermore, our start-to-end point guidance surpasses
98 their interactive refinement or reliance on 3D models in both efficiency and explainability. Some
99 additional studies [75, 88, 94, 14] investigated non-diffusion approaches to modeling embodied
100 world models.

101 2.2 End-to-End Vision-Language-Action Learning

102 End-to-end Vision-Language-Action (VLA) models constitute another major branch of embodied
103 learning, directly mapping visual and linguistic inputs to actions without explicit dynamics mod-
104eling. Recent advances in vision-language models (VLMs) [19, 48, 3, 49] have catalyzed the de-
105velopment of unified VLA systems that integrate perception, grounding, reasoning, and control
106within a single architecture—departing from traditional modular pipelines. At the core of most VLA
107systems lies the intuition that pre-trained vision-language models (VLMs) can serve as high-level
108semantic engines. Methods like RT-2 [11], OpenVLA [44], and RoboAgent [4] inject Internet-scale
109knowledge into robot policies by fine-tuning VLMs or conditioning them with low-level control to-
110kens. Based on their action modeling paradigms, VLA approaches can be broadly categorized into
111regression-based [53, 10, 58], generative/diffusion-based [7, 60, 83], and hybrid methods [59, 99].
112While these unified policies show strong performance on short-horizon or single-task scenarios, they
113often struggle with generalization and scalability—particularly across long horizons, domains, or em-
114bodyments. To address these limitations, recent efforts have introduced hierarchical structures into
115end-to-end VLA models, either by modularizing perception and policy [57, 82, 5], or by structuring
116execution into multi-stage skills [76, 70]. Though often labeled hierarchical, many still preserve end-
117to-end differentiability or latent consistency across stages. Such designs improve sample efficiency,
118handle long-horizon dependencies, and provide interpretable intermediate reasoning. Despite these

119 advances, existing VLA systems remain heavily dependent on in-domain demonstrations. Some
 120 mitigate this via cross-modal or cross-embodiment data [39, 67], or by scaling human video demon-
 121 strations [93, 84], yet diversity and coverage remain limited. Our work takes a complementary
 122 direction: we retain an end-to-end VLA formulation while introducing primitive-level hierarchy and
 123 spatially grounded guidance. This design enables zero-shot generalization across unseen tasks and
 124 settings, while maintaining modularity, interpretability, and execution efficiency.

125 **2.3 Datasets for Robotic Learning**

126 High-quality and diverse data is essential for generalizable robotic learning [92]. Existing large-
 127 scale datasets are typically collected via human teleoperation [9, 24, 42], scripted pipelines [22, 31],
 128 or expert policies [74], and have largely been consolidated in OpenX-Embodiment [22]. Despite its
 129 scale, this benchmark remains costly to scale and limited in semantic and morphological diversity.
 130 To enhance compositionality, RH20T-p [18] introduced primitive-level annotations atop RT-
 131 20T [28], but suffers from label noise and redundancy. Meanwhile, simulated benchmarks such
 132 as RLBench [40], CALVIN [65], and LIBERO [56] offer controllable environments, though often
 133 lack photorealism or task variety.
 134 Recent works [64, 30, 96, 91] explore augmenting human demos via simulation or generative ren-
 135 dering, but still rely on expert trajectories. In contrast, our framework directly generates large-scale,
 136 primitive-level demonstrations in simulation without human supervision or segmentation. Combined
 137 with pre-trained vision-language models, this enables efficient sim-to-real transfer with strong gen-
 138 eralization.

139 **3 Method**

140 Our approach is a modular and hierarchical system designed to bridge semantic-level planning and
 141 physically plausible action execution. It consists of three key components: a high-level semantic
 142 planner powered by a vision-language model (VLM), a primitive-conditioned video diffusion world
 143 model, and a zero-shot policy execution mechanism via trajectory projection and semi-closed-loop
 144 control. For a quick overview of our approach, please refer to Figure 1.

145 **3.1 Hierarchical Embodied World Modeling**

146 **3.1.1 Semantic Primitive Planning via Vision-Language Models**

147 Given a task instruction x^{text} and current observation frame x_0^{img} , we first employ a large vision-
 148 language model \mathcal{P} to perform semantic parsing and primitive planning:

$$\mathcal{A}_{1:N} = \mathcal{P}(x^{\text{text}}, x_0^{\text{img}}), \quad (1)$$

149 where $\mathcal{A}_{1:N}$ denotes the sequence of N atomic primitives (e.g., “pick up object A”, “place on B”).
 150 Each primitive is represented by an instruction a_k that is grounded locally in space and time.

151 For each primitive a_k , the planner identifies two pixel-space locations:

$$(s_k, g_k) = \mathcal{G}(a_k, x_0^{\text{img}}), \quad (2)$$

152 where $s_k \in \mathbb{R}^2$ and $g_k \in \mathbb{R}^2$ are the start and goal positions of the gripper in the image plane. These
 153 are Gaussian-blurred to produce heatmaps H_s and H_g , forming spatial guidance signals.

154 We leverage LoRA-based adaptation to enable reasoning and spatial grounding, following a Chain-
 155 of-LoRA [61]-style architecture. The structure is detailed in Figure 2.

156 **3.1.2 Primitive-Conditioned Video Diffusion Generation**

157 The video diffusion model \mathcal{D} is conditioned on the current frame x_0^{img} , the text description of the
 158 primitive a_k , and the spatial heatmaps H_s, H_g . The model learns to predict a rollout of T future
 159 frames:

$$x_{1:T}^{\text{img}} \sim \mathcal{D}(x_0^{\text{img}}, a_k, H_s, H_g). \quad (3)$$

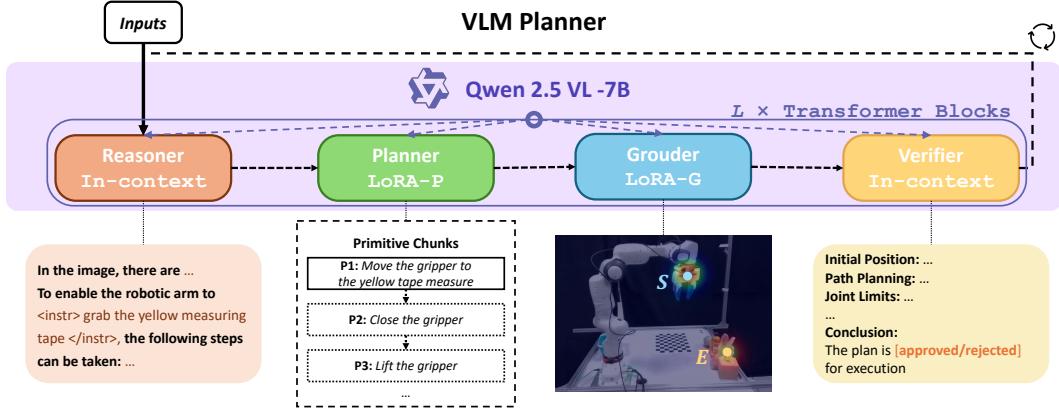


Figure 2: This workflow diagram delineates the architecture of our VLM planner, which is based on the Qwen 2.5 VL-7B model. The Reasoner generates descriptions and reasoning for the input to facilitate the functioning of subsequent components. The Planner (LoRA-P) and Grounder (LoRA-G) comprise two learnable LoRA modules: the Planner produces a series of primitive chunks, delivering modular instructions, while the Grounder provides specific descriptions of the end effector’s start and goal positions. The Verifier assesses the executability of the plan, indicating approved or rejected statuses to ensure precision in planning.

160 Unlike prior work that only observes cropped end-effector regions, \mathcal{D} is trained to observe the entire
 161 robot arm, allowing it to learn soft physical constraints such as reachability, base stability, and joint
 162 configuration limits.

163 The model is trained using a combination of pixel-level ℓ_2 reconstruction loss and perceptual simi-
 164 larity loss $\mathcal{L}_{\text{LPIPS}}$:

$$\mathcal{L}_{\text{vid}} = \sum_{t=1}^T \left\| x_t^{\text{img}} - \hat{x}_t^{\text{img}} \right\|_2^2 + \lambda \cdot \mathcal{L}_{\text{LPIPS}}(x_t^{\text{img}}, \hat{x}_t^{\text{img}}), \quad (4)$$

165 where \hat{x}_t^{img} denotes the ground truth frame at time t .

166 Appendix Figure 5 provides a detailed overview of the video generation module and the process of
 167 extracting 6-DoF end-effector trajectories. It should be explicitly noted that for primitives involving
 168 binary gripper actions (e.g., open or close), we bypass the video generation module and execute the
 169 action directly via symbolic control. This improves execution efficiency and enables a clean archi-
 170 tectural decoupling between discrete grasping commands and continuous 6-DoF motion planning,
 171 thereby enhancing modularity and interpretability.

172 3.2 Task-Agnostic Zero-Shot Policy and Execution

173 The generated future rollout $x_{1:T}^{\text{img}}$ for a primitive a_k lacks spatial information. To bridge the gap
 174 between pixel-space visual rollouts and real-world execution, we harness an off-the-shelf pose 6D
 175 estimation mechanism that requires only the generated RGB video as input.

176 Specifically, given a generated video sequence $V = \{I_1, I_2, \dots, I_T\}$ and a reference video $V_r =$
 177 $\{I_{r1}, I_{r2}, \dots, I_{rT'}\}$ of the robotic gripper, we estimate the gripper’s 6-DoF pose using a model-
 178 free, RGB-based pose estimator: $\mathbf{R}, \mathbf{T} = PE(V, V_r)$, where \mathbf{R} and \mathbf{T} denote the rotation and
 179 translation matrices, respectively. In our implementation, we use Gen6D [63] for its simplicity and
 180 generalization capability.

181 The pose information of the robotic gripper, particularly its 6-DoF object pose, utilizing the gener-
 182 ated video sequences directly for real-world robotic operations presents significant challenges.

183 After generating the future rollout $x_{1:T}^{\text{img}}$ for a primitive a_k , we extract a 6-DoF end-effector trajec-
 184 tory $\tau_k = \{p_1, \dots, p_T\}$ using 6-DoF pose estimation [63] and geometric transformation. These
 185 trajectories are then mapped into Cartesian control space and executed by the robot.

186 To project the estimated poses into real-world coordinates, we correct for the scale ambiguity inherent
 187 in monocular depth predictions. Given the depth map D at the initial frame and known camera
 188 intrinsics (x_c, y_c, f_x, f_y) , we compute the 3D location of any pixel (x, y) using:

$$X = (x - x_c) \cdot d/f_x, \quad Y = (y - y_c) \cdot d/f_y, \quad Z = d, \quad (5)$$

189 where d is the depth at pixel (x, y) . Since the generated video lacks absolute depth scale, we align the
 190 scale of the predicted trajectory using the ratio between the real-world depth d_0 (from the first frame
 191 of the depth camera) and the pixel-based depth d_0^{pixel} (from COLMAP or Gen6D). This yields a fixed
 192 correction factor $\lambda = \frac{d_0}{d_0^{\text{pixel}}}$, which is applied to all subsequent predicted translations: $\mathbf{T}_{\text{real}} = \lambda \cdot \mathbf{T}_{\text{pixel}}$.

193 The observed result of executing τ_k is then captured and passed back to the planner \mathcal{P} , enabling
 194 feedback-driven refinement in subsequent primitives. This forms a semi-closed-loop execution
 195 pipeline that mitigates error accumulation and enables task-aware adaptive control without end-to-
 196 end backpropagation.

197 Together, this design enables zero-shot generalization to unseen tasks, objects, and morphologies,
 198 using modular, interpretable, and resource-efficient world model components.

199 The process of mapping the generated video through the camera’s intrinsic and extrinsic parameters
 200 to obtain real-world coordinates is illustrated in Figure 3.

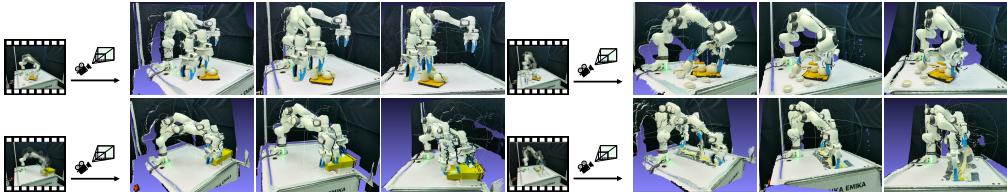


Figure 3: Mapping the generated video frames through the camera’s intrinsic and extrinsic parameters to visualize the corresponding 3D object poses in real-world coordinates.

201 3.3 Full-Arm, Multi-View, and Primitive-Aligned Data Collection Pipeline

202 Model generalization in embodied learning is fundamentally governed by the quality and structure
 203 of the available data [92]. Collecting large-scale and diverse data for embodied learning, however,
 204 remains notoriously labor-intensive. Prior real-world datasets often require years of human teleop-
 205 eration or scripted control to accumulate on the order of 100K-1M trajectories [9, 22, 43, 85, 23].
 206 Despite their size, such datasets typically exhibit narrow generalization, often struggling to extrapolate
 207 beyond minor variations in object pose or scene composition [67, 64]. Moreover, the need to
 208 manually segment tasks, align modalities, and annotate action labels imposes substantial overhead,
 209 limiting diversity and scalability.

210 To improve generalization in embodied agents, we advocate organizing and collecting data at the
 211 level of semantically grounded primitives. This structured design enables compositional generaliza-
 212 tion by allowing the model to cover a wide task space using a compact set of atomic operations. It
 213 also facilitates more focused representation learning, as the reduced temporal scope of each primitive
 214 helps the model capture salient object properties and interaction dynamics. Moreover, the semantic
 215 granularity simplifies annotation—coarse labels can be generated using vision-language models with
 216 minimal supervision. Finally, training the video diffusion model on short-horizon primitives sig-
 217 nificantly reduces computational overhead and eases optimization compared to long-horizon video
 218 prediction.

219 This rationale is further illustrated in appendix Figure 8, which presents a heuristic visualization of
 220 the impact of primitive structuring on generalization.

221 Specifically, we construct a primitive-centric dataset $\mathcal{D}_{\text{prim}}$ tailored for training our video diffusion
 222 world model \mathcal{W} . Each sample in $\mathcal{D}_{\text{prim}}$ consists of a tuple $(x_0^{\text{img}}, a_k, H_s, H_g, x_{1:T}^{\text{img}})$, where x_0^{img} is
 223 the initial observation, a_k is the primitive instruction, H_s and H_g are Gaussian heatmaps denoting
 224 gripper start and goal positions, and $x_{1:T}^{\text{img}}$ are the predicted future frames. Note that the primitive
 225 data can be either segmented during collecting or automatically segmented using motion heuristics
 226 and language-guided tagging tools, minimizing human supervision.

227 Another notable design choice in $\mathcal{D}_{\text{prim}}$ is that the full robot arm—including the base and joint
 228 structures—is always visible in the field of view, in contrast to cropped, end-effector-only frames
 229 common in prior datasets. This enables the diffusion model to internalize soft physical constraints
 230 such as reachability, joint limits, and spatial feasibility during prediction.

231 4 Experiments & Analysis

232 4.1 Dataset Description

233 Our dataset consists of 7,326 simulated and 11,465 real-world primitives collected using a 5-camera
 234 synchronized setup with Deoxys [104]. By segmenting long-horizon tasks via keyframe indices, we
 235 extracted on average 5.8 primitives per session, achieving up to 29 \times collection efficiency. As shown
 236 in Figure 4, cameras were arranged to maximize coverage of the workspace.

237 For annotation, 10% of the data was manually labeled, then used to fine-tune Qwen-VL-7B for
 238 auto-labeling the remainder, followed by light manual correction. To enhance visual generation
 239 quality, we mixed a small portion of diverse simulated data into the training set, including examples
 240 generated from RLBench, SayCan (PyBullet), and OpenVLA in LIBERO [56]. This hybrid strategy
 241 improves model generalization and dynamic realism, as confirmed in Section 4.4.

242 4.2 Baselines

243 We compare our framework against OpenVLA [44], a state-of-the-art end-to-end vision-language-
 244 action model that maps raw instructions and observations directly to actions. We evaluate OpenVLA
 245 under two settings: **(1) Zero-shot**: The pre-trained OpenVLA model is deployed directly on our
 246 benchmark tasks without any task-specific finetuning. This setting evaluates its raw generalization
 247 capability across unseen objects and scenes. **(2) Finetune**: The model is finetuned on 100 task-
 248 specific demonstrations per task using supervised behavior cloning, following the same protocol as
 249 in our system’s policy training. This setting represents an upper-bound of OpenVLAs performance
 250 under favorable conditions. We note that OpenVLA performs poorly in the zero-shot setting and
 251 fails completely on real-robot deployment, whereas our method maintains strong performance even
 252 without task-specific finetuning or retraining of the video model.

253 4.3 Performance

254 4.3.1 Overall Success Rate

Table 1: **Overall success rate on RLBench tasks.** We compare our method against existing baselines on 9 manipulation tasks from the RLBench benchmark [40], with success rates averaged over 100 episodes. Results for Image-BC [41], UniPi* [27], and 4DWM [100] are directly cited from the 4DWM paper [100]. Our method (bottom row) achieves the highest success rate on all tasks, with consistent gains in articulated and contact-rich scenarios.

Methods	close box	open drawer	open jar	open microwave	put knife	sweep to dustpan	lid off	weighing off	water plants
Image-BC	53	4	0	5	0	0	12	21	0
UniPi*	81	67	38	72	66	49	70	68	35
4DWM	88	80	44	70	70	56	73	62	41
Ours	93	84	43	78	72	63	67	58	56

255 Table 1 presents a comprehensive comparison of overall success rates across nine manipulation tasks
 256 from the RLBench benchmark . The performance of our method is assessed against several existing
 257 baselines, including Image-BC, UniPi*, and 4DWM, with success rates averaged over 100 episodes.

258 Our approach consistently outperforms all baseline methods across the various tasks. Notably, we
 259 achieve a success rate of 93% in the task of closing a box, which is the highest recorded. Similarly,
 260 our method attains a success rate of 84% in opening a drawer, further demonstrating its effectiveness
 261 in articulated scenarios. While the success rate for opening a jar is comparable to the baselines at
 262 43%, we maintain competitive performance in other tasks, such as opening a microwave (78%) and
 263 weighing plants (56%).

264 The results indicate that our method excels particularly in articulated and contact-rich tasks, where
265 the nuanced handling of objects is critical. These findings underscore the robustness and general-
266 ization capabilities of our approach, establishing it as a leading solution in the domain of robotic
267 manipulation.

268 **4.3.2 Planning**

Table 2: **Performance breakdown on three real-world tasks across planning, video generation, and primitive execution.** “Ours” denotes our method with frozen video diffusion; “OpenVLA” represents a strong end-to-end baseline. We additionally report OpenVLA’s zero-shot performance without task-specific fine-tuning.

Task	Stage	Metric	Ours	OpenVLA	OpenVLA (ZS)
Pick up cup	Planning	Primitive accuracy	18 / 20	N/A	N/A
	Video Generation	Frame realism (/total)	17 / 20	N/A	N/A
	Primitive Execution	Task success	16 / 20	12 / 20	0 / 20
Move cloth	Planning	Primitive accuracy	16 / 20	N/A	N/A
	Video Generation	Frame realism (/total)	15 / 20	N/A	N/A
	Primitive Execution	Task success	14 / 20	10 / 20	0 / 20
Fold cloth	Planning	Primitive accuracy	15 / 20	N/A	N/A
	Video Generation	Frame realism (/total)	14 / 20	N/A	N/A
	Primitive Execution	Task success	13 / 20	4 / 20	0 / 20

269 Our method performs a notable success rate, attributable to the effective decomposition of tasks into
270 manageable primitives. Each primitive is trained on a mixed dataset, which enhances its adaptability
271 and accuracy in real-world scenarios. Importantly, the test scenarios outlined in Table 2 were not
272 present in the training set, yet our method achieves a high primitive accuracy of 18 out of 20 for the
273 task of picking up a cup. Similar performance is observed across other tasks, with accuracy rates
274 of 16 out of 20 for moving cloth and 15 out of 20 for folding cloth. This robust planning capability
275 ensures that the robot can generate reliable sequences of actions, setting a strong foundation for
276 successful execution.

277 Notably, OpenVLA fails completely in the zero-shot setting, highlighting its reliance on task-specific
278 adaptation. In contrast, our method achieves strong generalization without retraining the video
279 model, benefiting from modularity and spatial grounding in its design.

280 **4.3.3 Primitive Execution**

281 The results further highlight the effectiveness of our approach. The integration of planning and
282 video generation significantly contributes to the overall task success rates. For instance, in the task
283 of picking up a cup, our method achieves a task success rate of 16 out of 20, compared to only
284 12 out of 20 for the OpenVLA baseline. Similar trends are observed in the other tasks, with our
285 method outperforming OpenVLA in both moving and folding cloth. Notably, the ability to succeed
286 in these tasks—despite their absence in the training set—demonstrates the robustness and generaliza-
287 tion capability of our method. This enhanced execution performance can be attributed to the careful
288 segmentation of primitives and the use of high-quality video generation, which collectively facilitate
289 precise and effective task completion.

290 **4.3.4 Long-Horizon Tasks Planning Performance**

291 The proposed methodology effectively addresses long-horizon tasks by decomposing them into a se-
292 ries of modular primitives, each designed to be plug-and-play. This hierarchical framework enables
293 each primitive to concentrate on short-range objectives, which aligns closely with the strengths of
294 Vision-Language Models (VLM). By constraining each component to operate within a limited con-
295 textual scope, we mitigate potential challenges such as contextual inconsistencies that may arise dur-
296 ing long-term planning. As illustrated in Figure 3, the process of mapping generated video frames to
297 real-world coordinates through the camera’s intrinsic and extrinsic parameters underpins this modu-
298 lar approach. This mapping facilitates precise pose estimations for each primitive, thereby ensuring
299 that the robotic system can execute tasks effectively without losing sight of the overarching long-

300 horizon goal. Consequently, our methodology leverages the inherent advantages of decomposing
301 complex tasks into manageable segments, enhancing both efficiency and reliability in execution.

302 4.4 Ablation Study

Table 3: **Ablation study on model components.** Task success rate is reported over 20 trials per setting. “Primitive planner ablations” isolate the effect of VLM-based spatial grounding; “VideoGen ablations” evaluate the impact of simulation-augmented training.

Ablation Group	Variant	Pick up cup	Move cloth	Fold cloth
Primitive Planner	Full model (with start/end prompts)	16 / 20	14 / 20	13 / 20
	w/o start-end prompt	12 / 20	10 / 20	7 / 20
	w/o primitive planner (direct instruction-to-action)	9 / 20	5 / 20	3 / 20
Video Generation	Full model (with sim + real data)	16 / 20	14 / 20	13 / 20
	Trained on real-only data	12 / 20	9 / 20	5 / 20

303 We conduct two ablation studies to assess the contribution of (1) the VLM-based primitive planner,
304 and (2) simulation-augmented training for video generation. As shown in Table 3, removing the
305 start-end spatial prompts from the planner leads to a noticeable decline in performance across all
306 tasks. This confirms that explicitly grounding primitives with spatial anchors significantly improves
307 downstream execution. A further removal of the entire primitive decomposition step—i.e., mapping
308 instructions directly to actions—results in severe performance degradation, especially on multi-step
309 or spatially entangled tasks like cloth folding. For video generation, removing simulated data from
310 the training corpus substantially weakens the model’s rollout quality, leading to lower execution
311 success rates. This suggests that simulation data offers critical diversity and coverage, allowing the
312 model to generalize better to real-world scenes with varied configurations and object appearances.

313 4.5 Efficiency

314 As shown in appendix Table 4, our model demonstrates significant advantages in both computational
315 efficiency and memory footprint. Compared to large-scale diffusion baselines such as Hunyuan I2V
316 and Wan 2.1, our method achieves up to **75x faster inference** and **67x lower VRAM usage**, making
317 it well suited for real-time robotic deployment. Despite its lightweight design, our method main-
318 tains competitive video quality and task relevance, validating its practicality in resource-constrained
319 scenarios.

320 5 Conclusion

321 We present a modular and interpretable framework for zero-shot embodied policy learning, com-
322 bining a vision-language planner with a primitive-conditioned video diffusion world model. By
323 grounding high-level instructions into spatially guided primitives and reusing a lightweight video
324 model across tasks, our system achieves strong performance on challenging real-world manipu-
325 lation benchmarks. Extensive experiments demonstrate superior generalization, data efficiency, and
326 runtime scalability. Our work highlights the potential of combining semantic reasoning with visual
327 dynamics modeling for scalable, reusable, and task-agnostic robotic intelligence.

328 6 Limitations

329 While our framework shows advancements in planning and decision-making for embodied agents,
330 several limitations warrant consideration:

331 **Challenges in High-Resolution Generation:** The framework struggles with high-resolution video
332 generation, impacting performance in tasks needing detailed visual fidelity.

333 **Dependency on Language Instructions:** It relies on vision-language models for deriving action
334 primitives, limiting adaptability to complex instructions not well-represented in the training data.

335 These limitations highlight areas for future research to enhance the framework’s applicability in
336 real-world scenarios.

337 **References**

- 338 [1] Alessandro Achille and Stefano Soatto. A separation principle for control in the age of deep
339 learning, 2017. URL <https://arxiv.org/abs/1711.03321>.
- 340 [2] Anurag Ajay, Seungwook Han, Yilun Du, Shuang Li, Abhi Gupta, Tommi Jaakkola, Josh
341 Tenenbaum, Leslie Kaelbling, Akash Srivastava, and Pulkit Agrawal. Compositional founda-
342 tion models for hierarchical planning. *Advances in Neural Information Processing Systems*,
343 36:22304–22325, 2023.
- 344 [3] Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang,
345 Daniel Salz, Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschannen, Emanuele
346 Bugliarello, et al. Paligemma: A versatile 3b vlm for transfer. *arXiv preprint arXiv:2407.07726*, 2024.
- 348 [4] Homanga Bharadhwaj, Jay Vakil, Mohit Sharma, Abhinav Gupta, Shubham Tulsiani, and
349 Vikash Kumar. Roboagent: Generalization and efficiency in robot manipulation via semantic
350 augmentations and action chunking, 2023.
- 351 [5] Johan Bjorck, Fernando Castañeda, Nikita Cherniadev, Xingye Da, Runyu Ding, Linxi Fan,
352 Yu Fang, Dieter Fox, Fengyuan Hu, Spencer Huang, et al. Gr00t n1: An open foundation
353 model for generalist humanoid robots. *arXiv preprint arXiv:2503.14734*, 2025.
- 354 [6] Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Nic-
355 colo Fusai, Lachy Groom, Karol Hausman, Brian Ichter, Szymon Jakubczak, Tim Jones,
356 Liyiming Ke, Sergey Levine, Adrian Li-Bell, Mohith Mothukuri, Suraj Nair, Karl Pertsch,
357 Lucy Xiaoyang Shi, James Tanner, Quan Vuong, Anna Walling, Haohuan Wang, and Ury
358 Zhilinsky. π_0 : A vision-language-action flow model for general robot control, 2024. URL
359 <https://arxiv.org/abs/2410.24164>.
- 360 [7] Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Nic-
361 colo Fusai, Lachy Groom, Karol Hausman, Brian Ichter, et al. π_0 : A vision-language-action
362 flow model for general robot control. *arXiv preprint arXiv:2410.24164*, 2024.
- 363 [8] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Do-
364 minik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, Varun Jampani, and
365 Robin Rombach. Stable video diffusion: Scaling latent video diffusion models to large
366 datasets, 2023. URL <https://arxiv.org/abs/2311.15127>.
- 367 [9] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea
368 Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al. Rt-1:
369 Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022.
- 370 [10] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof
371 Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, et al. Rt-2:
372 Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*, 2023.
- 374 [11] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof
375 Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, et al. Rt-2:
376 Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*, 2023.
- 378 [12] Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr,
379 Joe Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh.
380 Video generation models as world simulators. 2024. URL <https://openai.com/research/video-generation-models-as-world-simulators>.
- 382 [13] Jake Bruce, Michael D Dennis, Ashley Edwards, Jack Parker-Holder, Yuge Shi, Edward
383 Hughes, Matthew Lai, Aditi Mavalankar, Richie Steigerwald, Chris Apps, et al. Genie: Gener-
384 ative interactive environments. In *Forty-first International Conference on Machine Learning*,
385 2024.

- 386 [14] Chi-Lam Cheang, Guangzeng Chen, Ya Jing, Tao Kong, Hang Li, Yifeng Li, Yuxiao Liu,
 387 Hongtao Wu, Jiafeng Xu, Yichu Yang, Hanbo Zhang, and Minzhao Zhu. Gr-2: A genera-
 388 tive video-language-action model with web-scale knowledge for robot manipulation. *arXiv*
 389 *preprint arXiv:2410.06158*, 2024.
- 390 [15] Chang Chen, Yi-Fu Wu, Jaesik Yoon, and Sungjin Ahn. Transdreamer: Reinforcement learn-
 391 ing with transformer world models. *arXiv preprint arXiv:2202.09481*, 2022.
- 392 [16] Haoxin Chen, Menghan Xia, Yingqing He, Yong Zhang, Xiaodong Cun, Shaoshu Yang, Jinbo
 393 Xing, Yaofang Liu, Qifeng Chen, Xintao Wang, Chao Weng, and Ying Shan. Videocrafter1:
 394 Open diffusion models for high-quality video generation, 2023.
- 395 [17] Haoxin Chen, Yong Zhang, Xiaodong Cun, Menghan Xia, Xintao Wang, Chao Weng, and
 396 Ying Shan. Videocrafter2: Overcoming data limitations for high-quality video diffusion mod-
 397 els, 2024.
- 398 [18] Zeren Chen, Zhelun Shi, Xiaoya Lu, Lehan He, Sucheng Qian, Hao Shu Fang, Zhenfei Yin,
 399 Wanli Ouyang, Jing Shao, Yu Qiao, Cewu Lu, and Lu Sheng. Rh20t-p: A primitive-level
 400 robotic dataset towards composable generalization agents. *arXiv preprint arXiv: 2403.19622*,
 401 2024.
- 402 [19] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qing-
 403 long Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and
 404 aligning for generic visual-linguistic tasks. In *IEEE/CVF conference on computer vision and*
 405 *pattern recognition*, pages 24185–24198, 2024.
- 406 [20] Silvia Chiappa, Sébastien Racaniere, Daan Wierstra, and Shakir Mohamed. Recurrent envi-
 407 ronment simulators. *arXiv preprint arXiv:1704.02254*, 2017.
- 408 [21] Jaden Clark, Suvir Mirchandani, Dorsa Sadigh, and Suneel Belkhale. Action-free reasoning
 409 for policy generalization. In <https://arxiv.org/abs/2403.01823>, 2025.
- 410 [22] Open X-Embodiment Collaboration, Abby O'Neill, Abdul Rehman, Abhinav Gupta, Abhi-
 411 ram Maddukuri, Abhishek Gupta, Abhishek Padalkar, Abraham Lee, Acorn Pooley, Agrim
 412 Gupta, Ajay Mandlekar, Ajinkya Jain, Albert Tung, Alex Bewley, Alex Herzog, Alex Irpan,
 413 Alexander Khazatsky, Anant Rai, Anchit Gupta, Andrew Wang, Andrey Kolobov, Anikait
 414 Singh, Animesh Garg, Aniruddha Kembhavi, Annie Xie, Anthony Brohan, Antonin Raffin,
 415 Archit Sharma, Arefeh Yavary, Arhan Jain, Ashwin Balakrishna, Ayzaan Wahid, Ben Burgess-
 416 Limerick, Beomjoon Kim, Bernhard Schölkopf, Blake Wulfe, Brian Ichter, Cewu Lu, Charles
 417 Xu, Charlotte Le, Chelsea Finn, Chen Wang, Chenfeng Xu, Cheng Chi, Chenguang Huang,
 418 Christine Chan, Christopher Agia, Chuer Pan, Chuyuan Fu, Coline Devin, Danfei Xu, Daniel
 419 Morton, Danny Driess, Daphne Chen, Deepak Pathak, Dhruv Shah, Dieter Büchler, Dinesh
 420 Jayaraman, Dmitry Kalashnikov, Dorsa Sadigh, Edward Johns, Ethan Foster, Fangchen Liu,
 421 Federico Ceola, Fei Xia, Feiyu Zhao, Felipe Vieira Frujeri, Freek Stulp, Gaoyue Zhou, Gau-
 422 rav S. Sukhatme, Gautam Salhotra, Ge Yan, Gilbert Feng, Giulio Schiavi, Glen Berseth, Gre-
 423 gory Kahn, Guangwen Yang, Guanzhi Wang, Hao Su, Hao-Shu Fang, Haochen Shi, Henghui
 424 Bao, Heni Ben Amor, Henrik I Christensen, Hiroki Furuta, Homanga Bharadhwaj, Homer
 425 Walke, Hongjie Fang, Huy Ha, Igor Mordatch, Ilija Radosavovic, Isabel Leal, Jacky Liang,
 426 Jad Abou-Chakra, Jaehyung Kim, Jaimyn Drake, Jan Peters, Jan Schneider, Jasmine Hsu,
 427 Jay Vakil, Jeannette Bohg, Jeffrey Bingham, Jeffrey Wu, Jensen Gao, Jiaheng Hu, Jiajun Wu,
 428 Jialin Wu, Jiankai Sun, Jianlan Luo, Jiayuan Gu, Jie Tan, Jihoon Oh, Jimmy Wu, Jingpei Lu,
 429 Jingyun Yang, Jitendra Malik, João Silvério, Joey Hejna, Jonathan Booher, Jonathan Tomp-
 430 son, Jonathan Yang, Jordi Salvador, Joseph J. Lim, Junhyek Han, Kaiyuan Wang, Kanishka
 431 Rao, Karl Pertsch, Karol Hausman, Keegan Go, Keerthana Gopalakrishnan, Ken Goldberg,
 432 Kendra Byrne, Kenneth Oslund, Kento Kawaharazuka, Kevin Black, Kevin Lin, Kevin Zhang,
 433 Kiana Ehsani, Kiran Lekkala, Kirsty Ellis, Krishan Rana, Krishnan Srinivasan, Kuan Fang,
 434 Kunal Pratap Singh, Kuo-Hao Zeng, Kyle Hatch, Kyle Hsu, Laurent Itti, Lawrence Yunliang
 435 Chen, Lerrel Pinto, Li Fei-Fei, Liam Tan, Linxi Jim Fan, Lionel Ott, Lisa Lee, Luca Weihs,
 436 Magnum Chen, Marion Lepert, Marius Memmel, Masayoshi Tomizuka, Masha Itkina, Ma-
 437 teo Guaman Castro, Max Spero, Maximilian Du, Michael Ahn, Michael C. Yip, Mingtong
 438 Zhang, Mingyu Ding, Minho Heo, Mohan Kumar Srirama, Mohit Sharma, Moo Jin Kim,
 439 Naoaki Kanazawa, Nicklas Hansen, Nicolas Heess, Nikhil J Joshi, Niko Suenderhauf, Ning

Liu, Norman Di Palo, Nur Muhammad Mahi Shafullah, Oier Mees, Oliver Kroemer, Osbert Bastani, Pannag R Sanketi, Patrick Tree Miller, Patrick Yin, Paul Wohlhart, Peng Xu, Peter David Fagan, Peter Mitrano, Pierre Sermanet, Pieter Abbeel, Priya Sundaresan, Qiuyu Chen, Quan Vuong, Rafael Rafailov, Ran Tian, Ria Doshi, Roberto Mart'in-Mart'in, Rohan Baijal, Rosario Scalise, Rose Hendrix, Roy Lin, Runjia Qian, Ruohan Zhang, Russell Mendonca, Rutav Shah, Ryan Hoque, Ryan Julian, Samuel Bustamante, Sean Kirmani, Sergey Levine, Shan Lin, Sherry Moore, Shikhar Bahl, Shivin Dass, Shubham Sonawani, Shubham Tulsiani, Shuran Song, Sichun Xu, Siddhant Haldar, Siddharth Karamcheti, Simeon Adebola, Simon Guist, Soroush Nasiriany, Stefan Schaal, Stefan Welker, Stephen Tian, Subramanian Ramamoorthy, Sudeep Dasari, Suneel Belkhale, Sungjae Park, Suraj Nair, Suvir Mirchandani, Takayuki Osa, Tanmay Gupta, Tatsuya Harada, Tatsuya Matsushima, Ted Xiao, Thomas Kollar, Tianhe Yu, Tianli Ding, Todor Davchev, Tony Z. Zhao, Travis Armstrong, Trevor Darrell, Trinity Chung, Vidhi Jain, Vikash Kumar, Vincent Vanhoucke, Wei Zhan, Wenzuan Zhou, Wolfram Burgard, Xi Chen, Xiangyu Chen, Xiaolong Wang, Xinghao Zhu, Xinyang Geng, Xiyuan Liu, Xu Liangwei, Xuanlin Li, Yansong Pang, Yao Lu, Yecheng Jason Ma, Yejin Kim, Yevgen Chebotar, Yifan Zhou, Yifeng Zhu, Yilin Wu, Ying Xu, Yixuan Wang, Yonatan Bisk, Yongqiang Dou, Yoonyoung Cho, Youngwoon Lee, Yuchen Cui, Yue Cao, Yueh-Hua Wu, Yujin Tang, Yuke Zhu, Yunchu Zhang, Yunfan Jiang, Yunshuang Li, Yunzhu Li, Yusuke Iwasawa, Yutaka Matsuo, Zehan Ma, Zhuo Xu, Zichen Jeff Cui, Zichen Zhang, Zipeng Fu, and Zipeng Lin. Open X-Embodiment: Robotic learning datasets and RT-X models. <https://arxiv.org/abs/2310.08864>, 2023.

- [23] Murtaza Dalal, Min Liu, Walter Talbott, Chen Chen, Deepak Pathak, Jian Zhang, and Ruslan Salakhutdinov. Local policies enable zero-shot long-horizon manipulation. *International Conference of Robotics and Automation*, 2025.
- [24] Shivin Dass, Jullian Yapeter, Jesse Zhang, Jiahui Zhang, Karl Pertsch, Stefanos Nikolaidis, and Joseph J. Lim. Clvr jaco play dataset, 2023. URL https://github.com/clvrai/clvr_jaco_play_dataset.
- [25] Ming Ding, Wendi Zheng, Wenyi Hong, and Jie Tang. Cogview2: Faster and better text-to-image generation via hierarchical transformers, 2022. URL <https://arxiv.org/abs/2204.14217>.
- [26] Yilun Du and Leslie Kaelbling. Compositional generative modeling: A single model is not all you need. *arXiv preprint arXiv:2402.01103*, 2024.
- [27] 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*, 36:9156–9172, 2023.
- [28] Hao-Shu Fang, Hongjie Fang, Zhenyu Tang, Jirong Liu, Junbo Wang, Haoyi Zhu, and Cewu Lu. Rh20t: A robotic dataset for learning diverse skills in one-shot. In *RSS 2023 Workshop on Learning for Task and Motion Planning*, 2023.
- [29] Norm Ferns, Prakash Panangaden, and Doina Precup. Metrics for finite markov decision processes. 2004.
- [30] Caelan Garrett, Ajay Mandlekar, Bowen Wen, and Dieter Fox. Skillmimicgen: Automated demonstration generation for efficient skill learning and deployment, 2024. URL <https://arxiv.org/abs/2410.18907>.
- [31] Jiayuan Gu, Fanbo Xiang, Xuanlin Li, Zhan Ling, Xiqiang Liu, Tongzhou Mu, Yihe Tang, Stone Tao, Xinyue Wei, Yunchao Yao, Xiaodi Yuan, Pengwei Xie, Zhiao Huang, Rui Chen, and Hao Su. Maniskill2: A unified benchmark for generalizable manipulation skills. In *International Conference on Learning Representations*, 2023.
- [32] David Ha and Jürgen Schmidhuber. Recurrent world models facilitate policy evolution. *Advances in neural information processing systems*, 31, 2018.
- [33] Danijar Hafner, Timothy Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering atari with discrete world models. *arXiv preprint arXiv:2010.02193*, 2020.

- 492 [34] Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse con-
493 trol tasks through world models. *Nature*, pages 1–7, 2025.
- 494 [35] Yingqing He, Tianyu Yang, Yong Zhang, Ying Shan, and Qifeng Chen. Latent video diffu-
495 sion models for high-fidelity long video generation, 2023. URL <https://arxiv.org/abs/2211.13221>.
- 496 [36] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models, 2020.
498 URL <https://arxiv.org/abs/2006.11239>.
- 499 [37] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi,
500 and David J Fleet. Video diffusion models. In S. Koyejo, S. Mohamed,
501 A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Infor-
502 mation Processing Systems*, volume 35, pages 8633–8646. Curran Associates, Inc.,
503 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/39235c56aef13fb05a6adc95eb9d8d66-Paper-Conference.pdf.
- 504 [38] Yucheng Hu, Yanjiang Guo, Pengchao Wang, Xiaoyu Chen, Yen-Jen Wang, Jianke Zhang,
505 Koushil Sreenath, Chaochao Lu, and Jianyu Chen. Video prediction policy: A generalist
506 robot policy with predictive visual representations. *arXiv preprint arXiv:2412.14803*, 2024.
- 507 [39] Physical Intelligence, Kevin Black, Noah Brown, James Darpinian, Karan Dhabalia, Danny
508 Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, Manuel Y. Galliker, Dibya
509 Ghosh, Lachy Groom, Karol Hausman, Brian Ichter, Szymon Jakubczak, Tim Jones, Liyiming
510 Ke, Devin LeBlanc, Sergey Levine, Adrian Li-Bell, Mohith Mothukuri, Suraj Nair, Karl
511 Pertsch, Allen Z. Ren, Lucy Xiaoyang Shi, Laura Smith, Jost Tobias Springenberg, Kyle
512 Stachowicz, James Tanner, Quan Vuong, Homer Walke, Anna Walling, Haohuan Wang, Lili
513 Yu, and Ury Zhilinsky. $\pi_{0.5}$: a vision-language-action model with open-world generalization,
514 2025. URL <https://arxiv.org/abs/2504.16054>.
- 515 [40] Stephen James, Zicong Ma, David Rovick Arrojo, and Andrew J. Davison. Rlbench: The
516 robot learning benchmark & learning environment. *IEEE Robotics and Automation Letters*,
517 2020.
- 518 [41] Eric Jang, Alex Irpan, Mohi Khansari, Daniel Kappler, Frederik Ebert, Corey Lynch, Sergey
519 Levine, and Chelsea Finn. Bc-z: Zero-shot task generalization with robotic imitation learning.
520 In *Conference on Robot Learning*, pages 991–1002. PMLR, 2022.
- 521 [42] Alexander Khazatsky, Karl Pertsch, Suraj Nair, Ashwin Balakrishna, Sudeep Dasari, Sid-
522 dharth Karamcheti, Soroush Nasiriany, Mohan Kumar Srirama, Lawrence Yunliang Chen,
523 Kirsty Ellis, Peter David Fagan, Joey Hejna, Masha Itkina, Marion Lepert, Yecheng Jason
524 Ma, Patrick Tree Miller, Jimmy Wu, Suneel Belkhale, Shivin Dass, Huy Ha, Arhan Jain,
525 Abraham Lee, Youngwoon Lee, Marius Memmel, Sungjae Park, Ilija Radosavovic, Kaiyuan
526 Wang, Albert Zhan, Kevin Black, Cheng Chi, Kyle Beltran Hatch, Shan Lin, Jingpei Lu, Jean
527 Mercat, Abdul Rehman, Pannag R Sanketi, Archit Sharma, Cody Simpson, Quan Vuong,
528 Homer Rich Walke, Blake Wulfe, Ted Xiao, Jonathan Heewon Yang, Arefeh Yavary, Tony Z.
529 Zhao, Christopher Agia, Rohan Baijal, Mateo Guaman Castro, Daphne Chen, Qiuyu Chen,
530 Trinity Chung, Jaimyn Drake, Ethan Paul Foster, Jensen Gao, David Antonio Herrera, Minho
531 Heo, Kyle Hsu, Jiaheng Hu, Donovon Jackson, Charlotte Le, Yunshuang Li, Kevin Lin, Roy
532 Lin, Zehan Ma, Abhiram Maddukuri, Suvir Mirchandani, Daniel Morton, Tony Nguyen, Abi-
533 gail O’Neill, Rosario Scalise, Derick Seale, Victor Son, Stephen Tian, Emi Tran, Andrew E.
534 Wang, Yilin Wu, Annie Xie, Jingyun Yang, Patrick Yin, Yunchu Zhang, Osbert Bastani, Glen
535 Berseth, Jeannette Bohg, Ken Goldberg, Abhinav Gupta, Abhishek Gupta, Dinesh Jayaraman,
536 Joseph J Lim, Jitendra Malik, Roberto Martín-Martín, Subramanian Ramamoorthy, Dorsa
537 Sadigh, Shuran Song, Jiajun Wu, Michael C. Yip, Yuke Zhu, Thomas Kollar, Sergey Levine,
538 and Chelsea Finn. Droid: A large-scale in-the-wild robot manipulation dataset. 2024.
- 539 [43] Alexander Khazatsky, Karl Pertsch, Suraj Nair, Ashwin Balakrishna, Sudeep Dasari, Sid-
540 dharth Karamcheti, Soroush Nasiriany, Mohan Kumar Srirama, Lawrence Yunliang Chen,
541 Kirsty Ellis, et al. Droid: A large-scale in-the-wild robot manipulation dataset. *arXiv preprint
542 arXiv:2403.12945*, 2024.

- 544 [44] Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair,
 545 Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, et al. Openvla: An open-source
 546 vision-language-action model. *arXiv preprint arXiv:2406.09246*, 2024.
- 547 [45] Po-Chen Ko, Jiayuan Mao, Yilun Du, Shao-Hua Sun, and Joshua B Tenenbaum. Learning to
 548 Act from Actionless Videos through Dense Correspondences. *arXiv:2310.08576*, 2023.
- 549 [46] Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin
 550 Li, Bo Wu, Jianwei Zhang, et al. Hunyuanvideo: A systematic framework for large video
 551 generative models. *arXiv preprint arXiv:2412.03603*, 2024.
- 552 [47] Timothée Lesort, Natalia Díaz-Rodríguez, Jean-François Goudou, and David Filliat. State
 553 representation learning for control: An overview. *Neural Networks*, 108:379–392, 2018.
- 554 [48] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang,
 555 Peiyuan Zhang, Yanwei Li, Ziwei Liu, et al. Llava-onevision: Easy visual task transfer. *arXiv
 556 preprint arXiv:2408.03326*, 2024.
- 557 [49] Dingzhe Li, Yixiang Jin, Yuhao Sun, Hongze Yu, Jun Shi, Xiaoshuai Hao, Peng Hao, Huaping
 558 Liu, Fuchun Sun, Jianwei Zhang, et al. What foundation models can bring for robot learning
 559 in manipulation: A survey. *arXiv preprint arXiv:2404.18201*, 2024.
- 560 [50] Muyang Li, Ji Lin, Chenlin Meng, Stefano Ermon, Song Han, and Jun-Yan Zhu. Efficient
 561 spatially sparse inference for conditional gans and diffusion models. In S. Koyejo,
 562 S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural
 563 Information Processing Systems*, volume 35, pages 28858–28873. Curran Associates, Inc.,
 564 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/b9603de9e49d0838e53b6c9cf9d06556-Paper-Conference.pdf.
- 565 [51] Shuang Li, Yihuai Gao, Dorsa Sadigh, and Shuran Song. Unified video action model. *arXiv
 566 preprint arXiv:2503.00200*, 2025.
- 567 [52] Xin Li, Wenqing Chu, Ye Wu, Weihang Yuan, Fanglong Liu, Qi Zhang, Fu Li, Haocheng
 568 Feng, Errui Ding, and Jingdong Wang. Videogen: A reference-guided latent diffusion
 569 approach for high definition text-to-video generation, 2023. URL <https://arxiv.org/abs/2309.00398>.
- 570 [53] Xinghang Li, Minghuan Liu, Hanbo Zhang, Cunjun Yu, Jie Xu, Hongtao Wu, Chilam Cheang,
 571 Ya Jing, Weinan Zhang, Huaping Liu, et al. Vision-language foundation models as effective
 572 robot imitators. *arXiv preprint arXiv:2311.01378*, 2023.
- 573 [54] Junbang Liang, Ruoshi Liu, Ege Ozguroglu, Sruthi Sudhakar, Achal Dave, Pavel Tokmakov,
 574 Shuran Song, and Carl Vondrick. Dreamitate: Real-world visuomotor policy learning via
 575 video generation, 2024. URL <https://arxiv.org/abs/2406.16862>.
- 576 [55] Bin Lin, Yunyang Ge, Xinhua Cheng, Zongjian Li, Bin Zhu, Shaodong Wang, Xianyi He,
 577 Yang Ye, Shenghai Yuan, Liuhan Chen, et al. Open-sora plan: Open-source large video
 578 generation model. *arXiv preprint arXiv:2412.00131*, 2024.
- 579 [56] Bo Liu, Yifeng Zhu, Chongkai Gao, Yihao Feng, Qiang Liu, Yuke Zhu, and Peter Stone.
 580 Libero: Benchmarking knowledge transfer for lifelong robot learning. *arXiv preprint
 581 arXiv:2306.03310*, 2023.
- 582 [57] Jiaming Liu, Chenxuan Li, Guanqun Wang, Lily Lee, Kaichen Zhou, Sixiang Chen, Chuyan
 583 Xiong, Jiaxin Ge, Renrui Zhang, and Shanghang Zhang. Self-corrected multimodal large
 584 language model for end-to-end robot manipulation. *arXiv preprint arXiv:2405.17418*, 2024.
- 585 [58] Jiaming Liu, Mengzhen Liu, Zhenyu Wang, Pengju An, Xiaoqi Li, Kaichen Zhou, Senqiao
 586 Yang, Renrui Zhang, Yandong Guo, and Shanghang Zhang. Robomamba: Efficient vision-
 587 language-action model for robotic reasoning and manipulation. *Advances in Neural Informa-
 588 tion Processing Systems*, 37:40085–40110, 2024.
- 589 [59]

- 591 [59] Jiaming Liu, Hao Chen, Pengju An, Zhuoyang Liu, Renrui Zhang, Chenyang Gu, Xiaoqi
 592 Li, Ziyu Guo, Sixiang Chen, Mengzhen Liu, et al. Hybridvla: Collaborative diffusion and
 593 autoregression in a unified vision-language-action model. *arXiv preprint arXiv:2503.10631*,
 594 2025.
- 595 [60] Songming Liu, Lingxuan Wu, Bangguo Li, Hengkai Tan, Huayu Chen, Zhengyi Wang, Ke Xu,
 596 Hang Su, and Jun Zhu. Rdt-1b: a diffusion foundation model for bimanual manipulation.
 597 *arXiv preprint arXiv:2410.07864*, 2024.
- 598 [61] Ye Liu, Kevin Qinghong Lin, Chang Wen Chen, and Mike Zheng Shou. Videomind: A
 599 chain-of-lora agent for long video reasoning. *arXiv preprint arXiv:2503.13444*, 2025.
- 600 [62] Yixin Liu, Kai Zhang, Yuan Li, Zhiling Yan, Chujie Gao, Ruoxi Chen, Zhengqing Yuan,
 601 Yue Huang, Hanchi Sun, Jianfeng Gao, et al. Sora: A review on background, technology,
 602 limitations, and opportunities of large vision models. *arXiv preprint arXiv:2402.17177*, 2024.
- 603 [63] Yuan Liu, Yilin Wen, Sida Peng, Cheng Lin, Xiaoxiao Long, Taku Komura, and Wenping
 604 Wang. Gen6d: Generalizable model-free 6-dof object pose estimation from rgb images. In
 605 *Proceedings of the European Conference on Computer Vision*, pages 298–315, 2022.
- 606 [64] Ajay Mandlekar, Soroush Nasiriany, Bowen Wen, Iretiayo Akinola, Yashraj Narang, Linxi
 607 Fan, Yuke Zhu, and Dieter Fox. Mimicgen: A data generation system for scalable robot
 608 learning using human demonstrations, 2023. URL <https://arxiv.org/abs/2310.17596>.
- 609 [65] Oier Mees, Lukas Hermann, Erick Rosete-Beas, and Wolfram Burgard. Calvin: A benchmark
 610 for language-conditioned policy learning for long-horizon robot manipulation tasks. *IEEE*
 611 *Robotics and Automation Letters (RA-L)*, 7(3):7327–7334, 2022.
- 612 [66] Vincent Micheli, Eloi Alonso, and François Fleuret. Transformers are sample-efficient world
 613 models. *arXiv preprint arXiv:2209.00588*, 2022.
- 614 [67] NVIDIA, :, Johan Bjorck, Fernando Castañeda, Nikita Cherniadev, Xingye Da, Runyu Ding,
 615 Linxi Jim Fan, Yu Fang, Dieter Fox, Fengyuan Hu, Spencer Huang, Joel Jang, Zhenyu Jiang,
 616 Jan Kautz, Kaushil Kundalia, Lawrence Lao, Zhiqi Li, Zongyu Lin, Kevin Lin, Guilin Liu,
 617 Edith Llontop, Loïc Magne, Ajay Mandlekar, Avnish Narayan, Soroush Nasiriany, Scott
 618 Reed, You Liang Tan, Guanzhi Wang, Zu Wang, Jing Wang, Qi Wang, Jiannan Xiang, Yuqi
 619 Xie, Yinzen Xu, Zhenjia Xu, Seonghyeon Ye, Zhiding Yu, Ao Zhang, Hao Zhang, Yizhou
 620 Zhao, Ruijie Zheng, and Yuke Zhu. Gr0ot n1: An open foundation model for generalist
 621 humanoid robots, 2025. URL <https://arxiv.org/abs/2503.14734>.
- 622 [68] William Peebles and Saining Xie. Scalable diffusion models with transformers, 2023. URL
 623 <https://arxiv.org/abs/2212.09748>.
- 624 [69] Karl Pertsch, Kyle Stachowicz, Brian Ichter, Danny Driess, Suraj Nair, Quan Vuong, Oier
 625 Mees, Chelsela Finn, and Sergey Levine. Fast: Efficient action tokenization for vision-
 626 language-action models. *arXiv preprint arXiv:2501.09747*, 2025.
- 627 [70] Physical Intelligence. π 0.5: A vision-language-action model with open-world generalization.
 628 <https://www.physicalintelligence.company/blog/pi05>, 2025. Accessed: 2025-04-
 629 25.
- 630 [71] Adam Polyak, Amit Zohar, Andrew Brown, Andros Tjandra, Animesh Sinha, Ann Lee,
 631 Apoorv Vyas, Bowen Shi, Chih-Yao Ma, Ching-Yao Chuang, David Yan, Dhruv Choudhary,
 632 Dingkang Wang, Geet Sethi, Guan Pang, Haoyu Ma, Ishan Misra, Ji Hou, Jialiang
 633 Wang, Kiran Jagadeesh, Kunpeng Li, Luxin Zhang, Mannat Singh, Mary Williamson, Matt
 634 Le, Matthew Yu, Mitesh Kumar Singh, Peizhao Zhang, Peter Vajda, Quentin Duval, Ro-
 635 hit Girdhar, Roshan Sumbaly, Sai Saketh Rambhatla, Sam Tsai, Samaneh Azadi, Samyak
 636 Datta, Sanyuan Chen, Sean Bell, Sharadh Ramaswamy, Shelly Sheynin, Siddharth Bhat-
 637 tacharya, Simran Motwani, Tao Xu, Tianhe Li, Tingbo Hou, Wei-Ning Hsu, Xi Yin, Xiaoliang
 638 Dai, Yaniv Taigman, Yaqiao Luo, Yen-Cheng Liu, Yi-Chiao Wu, Yue Zhao, Yuval Kirstain,
 639 Zecheng He, Zijian He, Albert Pumarola, Ali Thabet, Artsiom Sanakoyeu, Arun Mallya, Bais-
 640 han Guo, Boris Araya, Breena Kerr, Carleigh Wood, Ce Liu, Cen Peng, Dmitry Vengertsev,

- 641 Edgar Schonfeld, Elliot Blanchard, Felix Juefei-Xu, Fraylie Nord, Jeff Liang, John Hoff-
 642 man, Jonas Kohler, Kaolin Fire, Karthik Sivakumar, Lawrence Chen, Licheng Yu, Luya Gao,
 643 Markos Georgopoulos, Rashel Moritz, Sara K. Sampson, Shikai Li, Simone Parmeggiani,
 644 Steve Fine, Tara Fowler, Vladan Petrovic, and Yuming Du. Movie gen: A cast of media
 645 foundation models, 2025. URL <https://arxiv.org/abs/2410.13720>.
- 646 [72] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical
 647 text-conditional image generation with clip latents, 2022. URL <https://arxiv.org/abs/2204.06125>.
- 649 [73] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer.
 650 High-resolution image synthesis with latent diffusion models, 2022. URL <https://arxiv.org/abs/2112.10752>.
- 652 [74] Giulio Schiavi, Paula Wulkop, Giuseppe Rizzi, Lionel Ott, Roland Siegwart, and Jen Jen
 653 Chung. Learning agent-aware affordances for closed-loop interaction with articulated objects.
 654 In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5916–
 655 5922. IEEE, 2023.
- 656 [75] Younggyo Seo, Danijar Hafner, Hao Liu, Fangchen Liu, Stephen James, Kimin Lee, and
 657 Pieter Abbeel. Masked world models for visual control. In *Conference on Robot Learning*,
 658 pages 1332–1344. PMLR, 2023.
- 659 [76] Lucy Xiaoyang Shi, Brian Ichter, Michael Equi, Liyiming Ke, Karl Pertsch, Quan Vuong,
 660 James Tanner, Anna Walling, Haohuan Wang, Niccolo Fusai, et al. Hi robot: Open-
 661 ended instruction following with hierarchical vision-language-action models. *arXiv preprint arXiv:2502.19417*, 2025.
- 663 [77] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan
 664 Hu, Harry Yang, Oron Ashual, Oran Gafni, Devi Parikh, Sonal Gupta, and Yaniv Taig-
 665 man. Make-a-video: Text-to-video generation without text-video data, 2022. URL <https://arxiv.org/abs/2209.14792>.
- 667 [78] Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep
 668 unsupervised learning using nonequilibrium thermodynamics, 2015. URL <https://arxiv.org/abs/1503.03585>.
- 670 [79] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models, 2022.
 671 URL <https://arxiv.org/abs/2010.02502>.
- 672 [80] Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon,
 673 and Ben Poole. Score-based generative modeling through stochastic differential equations,
 674 2021. URL <https://arxiv.org/abs/2011.13456>.
- 675 [81] Aether Team, Haoyi Zhu, Yifan Wang, Jianjun Zhou, Wenzheng Chang, Yang Zhou, Zizun Li,
 676 Junyi Chen, Chunhua Shen, Jiangmiao Pang, et al. Aether: Geometric-aware unified world
 677 modeling. *arXiv preprint arXiv:2503.18945*, 2025.
- 682 [82] Gemini Robotics Team, Saminda Abeyruwan, Joshua Ainslie, Jean-Baptiste Alayrac,
 683 Montserrat Gonzalez Arenas, Travis Armstrong, Ashwin Balakrishna, Robert Baruch, Maria
 684 Bauza, Michiel Blokzijl, et al. Gemini robotics: Bringing ai into the physical world. *arXiv preprint arXiv:2503.20020*, 2025.
- 685 [83] Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees,
 686 Sudeep Dasari, Joey Hejna, Tobias Kreiman, Charles Xu, et al. Octo: An open-source generalist
 687 robot policy. *arXiv preprint arXiv:2405.12213*, 2024.
- 688 [84] Yang Tian, Sizhe Yang, Jia Zeng, Ping Wang, Dahua Lin, Hao Dong, and Jiangmiao Pang.
 689 Predictive inverse dynamics models are scalable learners for robotic manipulation, 2024. URL
 690 <https://arxiv.org/abs/2412.15109>.

- 688 [85] Homer Walke, Kevin Black, Abraham Lee, Moo Jin Kim, Max Du, Chongyi Zheng, Tony
 689 Zhao, Philippe Hansen-Estruch, Quan Vuong, Andre He, Vivek Myers, Kuan Fang, Chelsea
 690 Finn, and Sergey Levine. Bridgedata v2: A dataset for robot learning at scale. In *Conference*
 691 *on Robot Learning (CoRL)*, 2023.
- 692 [86] Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu, Haim-
 693 ing Zhao, Jianxiao Yang, Jianyuan Zeng, Jiayu Wang, Jingfeng Zhang, Jingren Zhou, Jinkai
 694 Wang, Jixuan Chen, Kai Zhu, Kang Zhao, Keyu Yan, Lianghua Huang, Mengyang Feng,
 695 Ningyi Zhang, Pandeng Li, Pingyu Wu, Ruihang Chu, Ruili Feng, Shiwei Zhang, Siyang
 696 Sun, Tao Fang, Tianxing Wang, Tianyi Gui, Tingyu Weng, Tong Shen, Wei Lin, Wei Wang,
 697 Wei Wang, Wenmeng Zhou, Wente Wang, Wenting Shen, Wenyuan Yu, Xianzhong Shi, Xi-
 698 aoming Huang, Xin Xu, Yan Kou, Yangyu Lv, Yifei Li, Yijing Liu, Yiming Wang, Yingya
 699 Zhang, Yitong Huang, Yong Li, You Wu, Yu Liu, Yulin Pan, Yun Zheng, Yuntao Hong, Yu-
 700 peng Shi, Yutong Feng, Zeyinzi Jiang, Zhen Han, Zhi-Fan Wu, and Ziyu Liu. Wan: Open and
 701 advanced large-scale video generative models. *arXiv preprint arXiv:2503.20314*, 2025.
- 702 [87] Jiuniu Wang, Hangjie Yuan, Dayou Chen, Yingya Zhang, Xiang Wang, and Shiwei Zhang.
 703 Modelscope text-to-video technical report, 2023. URL <https://arxiv.org/abs/2308.06571>.
- 705 [88] Hongtao Wu, Ya Jing, Chilam Cheang, Guangzeng Chen, Jiafeng Xu, Xinghang Li,
 706 Minghuan Liu, Hang Li, and Tao Kong. Unleashing large-scale video generative pre-training
 707 for visual robot manipulation. In *International Conference on Learning Representations*,
 708 2024.
- 709 [89] Jiannan Xiang, Guangyi Liu, Yi Gu, Qiyue Gao, Yuting Ning, Yuheng Zha, Zeyu Feng, Tian-
 710 hua Tao, Shibo Hao, Yemin Shi, et al. Pandora: Towards general world model with natural
 711 language actions and video states. *arXiv preprint arXiv:2406.09455*, 2024.
- 712 [90] Jinbo Xing, Menghan Xia, Yong Zhang, Haoxin Chen, Wangbo Yu, Hanyuan Liu, Xintao
 713 Wang, Tien-Tsin Wong, and Ying Shan. Dynamicrafter: Animating open-domain images
 714 with video diffusion priors. *arXiv preprint arXiv:2310.12190*, 2023.
- 715 [91] Zhengrong Xue, Shuying Deng, Zhenyang Chen, Yixuan Wang, Zhecheng Yuan, and Huazhe
 716 Xu. Demogen: Synthetic demonstration generation for data-efficient visuomotor policy learn-
 717 ing, 2025. URL <https://arxiv.org/abs/2502.16932>.
- 718 [92] Zhengrong Xue, Shuying Deng, Zhenyang Chen, Yixuan Wang, Zhecheng Yuan, and Huazhe
 719 Xu. Demogen: Synthetic demonstration generation for data-efficient visuomotor policy learn-
 720 ing. *arXiv preprint arXiv:2502.16932*, 2025.
- 721 [93] Jiange Yang, Bei Liu, Jianlong Fu, Bocheng Pan, Gangshan Wu, and Limin Wang. Spatiotem-
 722 poral predictive pre-training for robotic motor control. *arXiv preprint arXiv:2403.05304*,
 723 2024.
- 724 [94] Jiange Yang, Bei Liu, Jianlong Fu, Bocheng Pan, Gangshan Wu, and Limin Wang. Spa-
 725 tiotemporal predictive pre-training for robotic motor control, 2024. URL <https://arxiv.org/abs/2403.05304>.
- 727 [95] Mengjiao Yang, Yilun Du, Kamyar Ghasemipour, Jonathan Tompson, Dale Schuurmans, and
 728 Pieter Abbeel. Learning interactive real-world simulators. *arXiv preprint arXiv:2310.06114*,
 729 2023.
- 730 [96] Sizhe Yang, Wenye Yu, Jia Zeng, Jun Lv, Kerui Ren, Cewu Lu, Dahua Lin, and Jiangmiao
 731 Pang. Novel demonstration generation with gaussian splatting enables robust one-shot ma-
 732 nipulation, 2025. URL <https://arxiv.org/abs/2504.13175>.
- 733 [97] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming
 734 Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion
 735 models with an expert transformer. *arXiv preprint arXiv:2408.06072*, 2024.

- 736 [98] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming
 737 Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, Da Yin, Yuxuan Zhang, Weihan Wang,
 738 Yean Cheng, Bin Xu, Xiaotao Gu, Yuxiao Dong, and Jie Tang. Cogvideox: Text-to-video
 739 diffusion models with an expert transformer, 2025. URL <https://arxiv.org/abs/2408.06072>.
- 741 [99] Seonghyeon Ye, Joel Jang, Byeongguk Jeon, Sejune Joo, Jianwei Yang, Baolin Peng, Ajay
 742 Mandlekar, Reuben Tan, Yu-Wei Chao, Bill Yuchen Lin, et al. Latent action pretraining from
 743 videos. *arXiv preprint arXiv:2410.11758*, 2024.
- 744 [100] Haoyu Zhen, Qiao Sun, Hongxin Zhang, Junyan Li, Siyuan Zhou, Yilun Du, and Chuang
 745 Gan. Tesseract: Learning 4d embodied world models. 2025. URL <https://arxiv.org/abs/2504.20995>.
- 747 [101] Zangwei Zheng, Xiangyu Peng, Tianji Yang, Chenhui Shen, Shenggui Li, Hongxin Liu,
 748 Yukun Zhou, Tianyi Li, and Yang You. Open-sora: Democratizing efficient video produc-
 749 tion for all, 2024. URL <https://arxiv.org/abs/2412.20404>.
- 750 [102] Siyuan Zhou, Yilun Du, Jaben Chen, Yandong Li, Dit-Yan Yeung, and Chuang Gan. Ro-
 751 bodreamer: Learning compositional world models for robot imagination, 2024. URL
 752 <https://arxiv.org/abs/2404.12377>.
- 753 [103] Chunming Zhu, Raymond Yu, Siyuan Feng, Benjamin Burchfiel, Paarth Shah, and Abhishek
 754 Gupta. Unified world models: Coupling video and action diffusion for pretraining on large
 755 robotic datasets. In *Proceedings of Robotics: Science and Systems (RSS)*, 2025.
- 756 [104] Yifeng Zhu, Abhishek Joshi, Peter Stone, and Yuke Zhu. Viola: Imitation learning for vision-
 757 based manipulation with object proposal priors. *6th Annual Conference on Robot Learning
 758 (CoRL)*, 2022.

759 A Experimental Environment

760 Our experimental setup is illustrated in Figure 4. The environment consists of a workstation, a FR3
 761 (Franka Research 3) robotic arm, and five cameras. The specific configuration is as follows:

762 **Workstation:** The workstation serves as a standard experimental platform, providing a stable area
 763 to support various operational tasks.

764 **FR3 Robotic Arm:** The FR3 robotic arm is a high-precision industrial robot equipped with flexible
 765 movement capabilities and high repeatability. This robotic arm is responsible for executing various
 766 tasks and can interact in real-time with the collected visual data.

767 **Camera System:** - **Two Realsense Cameras:** These cameras are mounted on the wrist of the
 768 robotic arm and on the shelf of the workstation, respectively. They capture real-time depth informa-
 769 tion and the arm’s movements, enabling a comprehensive understanding of the spatial context and
 770 enhancing the accuracy and reliability of task execution. - **Three Femto Bolt Cameras:** Positioned
 771 around the workstation, these cameras are used to capture dynamic processes at high frame rates.
 772 They provide additional visual perspectives, ensuring that comprehensive image data is collected
 773 during complex tasks for subsequent analysis and processing.

774 In summary, the design of this experimental environment aims to provide comprehensive visual
 775 support for robotic operations, facilitating the research and development of complex tasks. Through
 776 the collaborative operation of a multi-camera system, we can obtain high-quality visual data, thereby
 777 enhancing the performance of the robot in practical applications.

778 B Heuristic Intuition for Primitive-Level Data Organization Enabling 779 Compositional Generalization under Sparse Embodied Data

780 Figure 8 presents a heuristic Intuition for Primitive-Level Data Organization Enabling Composi-
 781 tional Generalization under Sparse Embodied Data.

782 **C An Illustration for Primitive-Level Task Execution**

783 Figure 5 illustrates a primitive-level task execution for Pick up the yellow tape measure. Motion
784 primitives are executed via diffusion-based visual rollout and trajectory projection, while discrete
785 gripper actions are handled directly through symbolic execution.

786 **D Examples of Training Video Clips**

787 In this study, we utilize both real-world and simulated video clips capturing the operation of robotic
788 arms to train our model. The real-world data consists of videos collected from various robotic arm
789 tasks using a Franka Emika arm, remotely operated through a Space Mouse interface. These tasks
790 are performed in controlled environments, showcasing different motions and interactions with ob-
791 jects. These videos provide valuable information about how the arm operates in realistic conditions,
792 with varying lighting, camera angles, and object placements.

793 The simulated video clips, on the other hand, are generated from physics-based environments such
794 as **RLBench** and **LIBERO**. These simulations replicate the robotic arm's movements in virtual
795 settings, allowing for the generation of large quantities of data under controlled conditions. This en-
796 ables the model to learn from a diverse range of scenarios that might be difficult or time-consuming
797 to capture in real life.

798 Examples of both real-world and simulated video clips are provided in Figure 6. These images
799 illustrate the types of data used to train the model, offering a glimpse into the varied nature of the
800 training set.

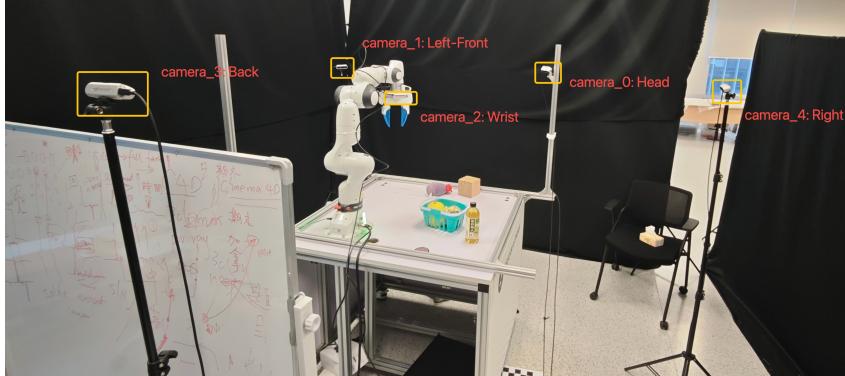


Figure 4: The workstation for data collection and real-robotic evaluation.

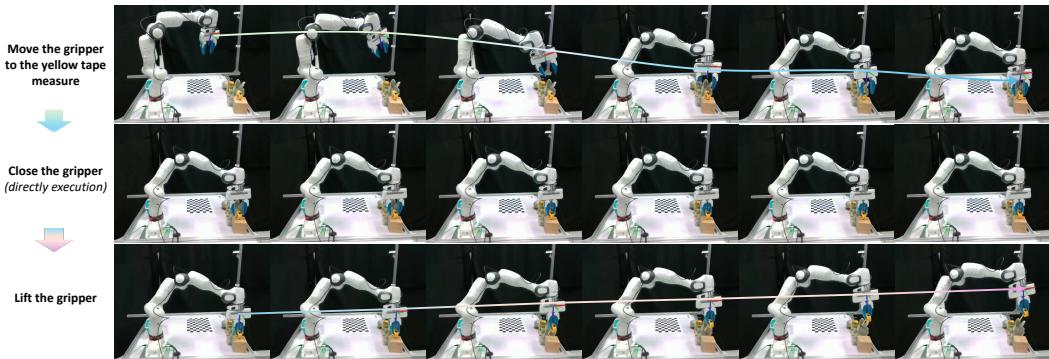


Figure 5: Illustration of primitive-level task execution for Pick up the yellow tape measure. Motion primitives are executed via diffusion-based visual rollout and trajectory projection, while discrete gripper actions are handled directly through symbolic execution.



Figure 6: Sample frames from the training dataset, which consists of three components: real-world data collected using a Franka Emika robotic arm and Femto Bolt cameras, and two simulated datasets obtained from RLBench and LIBERO.

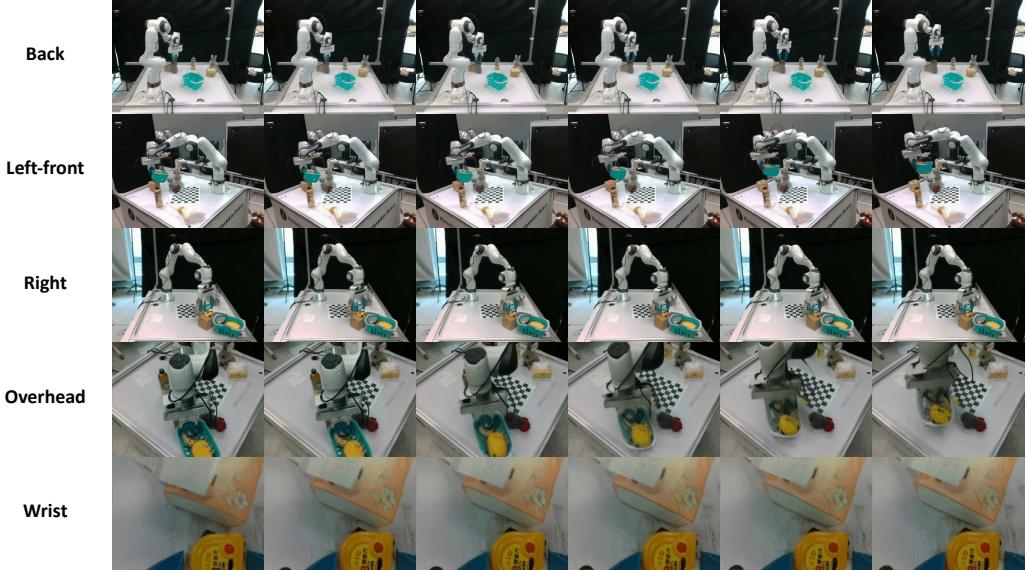


Figure 7: Sample frames from the generated video clips showing the Franka robotic arm performing manipulation tasks from five distinct viewpoints: **Back**, **Left-front**, **Right**, **Overhead**, and **Wrist**. Each viewpoint provides a unique perspective on the arm’s movements, enhancing the model’s ability to learn complex manipulation behaviors from diverse angles.

801 E Examples of Generated Video Clips

802 In addition to using recorded video data, we also generate synthetic video clips to enrich the training
803 dataset and enhance the models generalization capability. These generated clips simulate the Franka
804 robotic arm performing various manipulation tasks from multiple camera perspectives.

805 Specifically, each video sequence is rendered from five distinct viewpoints: **Back**, **Left-front**, **Right**,
806 **Overhead**, and **Wrist**. These perspectives are selected to comprehensively capture the arm’s kinematics,
807 end-effector trajectories, and object interactions from both global and local contexts. The
808 Back and Overhead views provide an overall understanding of the workspace and arm configuration,
809 while Left-front and Right offer lateral angles that help reveal occluded motions. The Wrist
810 view, positioned close to the end-effector, offers detailed observation of grasping and manipulation
811 actions.

812 Representative frames from these generated clips are shown in Figure 7, illustrating how each view-
813 point contributes to a multi-faceted understanding of the robotic operation.

814 **F Efficiency Comparisons**

815 Table 4 compares the video generation efficiency.

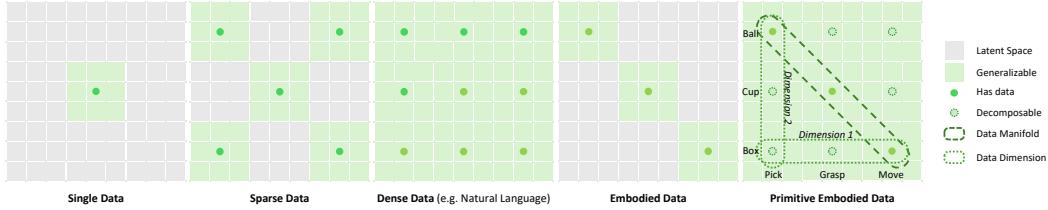


Figure 8: While densely distributed data can enable broad generalization, embodied data often suffer from sparsity [26, 92]. The rightmost schematic highlights how organizing embodied data at the primitive level—along orthogonal dimensions such as action and object—supports compositional generalization even under limited data availability.

Table 4: **Efficiency comparison of video generation models.** We compare video generation speed, memory footprint, and runtime on A100 GPUs. Ours achieves the best speed-VRAM trade-off under realistic deployment conditions.

Model	Resolution	VRAM (A100)	Time / Frames	FPS
Hunyuan I2V	480p	6079 GB	50 min / 81 frames (local)	0.027
4DWM (CogVideoX1.5-5B-I2)	480p	20 GB	18m20s / 49 frames	0.045
Wan 2.1 I2V (14B)	720p	76.7 GB	2715s / 81 frames	0.03
Ours	480p	11 GB	16s / 32 frames	2.0

816 **NeurIPS Paper Checklist**

817 **1. Claims**

818 Question: Do the main claims made in the abstract and introduction accurately reflect the
819 paper's contributions and scope?

820 Answer: [Yes]

821 Justification: The main claims made in the abstract and introduction accurately reflect the
822 paper's contributions and scope.

823 Guidelines:

- 824 • The answer NA means that the abstract and introduction do not include the claims
825 made in the paper.
- 826 • The abstract and/or introduction should clearly state the claims made, including the
827 contributions made in the paper and important assumptions and limitations. A No or
828 NA answer to this question will not be perceived well by the reviewers.
- 829 • The claims made should match theoretical and experimental results, and reflect how
830 much the results can be expected to generalize to other settings.
- 831 • It is fine to include aspirational goals as motivation as long as it is clear that these
832 goals are not attained by the paper.

833 **2. Limitations**

834 Question: Does the paper discuss the limitations of the work performed by the authors?

835 Answer: [Yes]

836 Justification: The paper discuss the limitations of the work in section "Limitations".

837 Guidelines:

- 838 • The answer NA means that the paper has no limitation while the answer No means
839 that the paper has limitations, but those are not discussed in the paper.
- 840 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 841 • The paper should point out any strong assumptions and how robust the results are to
842 violations of these assumptions (e.g., independence assumptions, noiseless settings,
843 model well-specification, asymptotic approximations only holding locally). The au-
844 thors should reflect on how these assumptions might be violated in practice and what
845 the implications would be.
- 846 • The authors should reflect on the scope of the claims made, e.g., if the approach was
847 only tested on a few datasets or with a few runs. In general, empirical results often
848 depend on implicit assumptions, which should be articulated.
- 849 • The authors should reflect on the factors that influence the performance of the ap-
850 proach. For example, a facial recognition algorithm may perform poorly when image
851 resolution is low or images are taken in low lighting. Or a speech-to-text system might
852 not be used reliably to provide closed captions for online lectures because it fails to
853 handle technical jargon.
- 854 • The authors should discuss the computational efficiency of the proposed algorithms
855 and how they scale with dataset size.
- 856 • If applicable, the authors should discuss possible limitations of their approach to ad-
857 dress problems of privacy and fairness.
- 858 • While the authors might fear that complete honesty about limitations might be used by
859 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
860 limitations that aren't acknowledged in the paper. The authors should use their best
861 judgment and recognize that individual actions in favor of transparency play an impor-
862 tant role in developing norms that preserve the integrity of the community. Reviewers
863 will be specifically instructed to not penalize honesty concerning limitations.

864 **3. Theory assumptions and proofs**

865 Question: For each theoretical result, does the paper provide the full set of assumptions and
866 a complete (and correct) proof?

867 Answer: [Yes]

868 Justification: A theoretical proof of the methods used is provided in the "Method" section.
869 Guidelines:

- 870 • The answer NA means that the paper does not include theoretical results.
- 871 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
872 referenced.
- 873 • All assumptions should be clearly stated or referenced in the statement of any theo-
874 rems.
- 875 • The proofs can either appear in the main paper or the supplemental material, but if
876 they appear in the supplemental material, the authors are encouraged to provide a
877 short proof sketch to provide intuition.
- 878 • Inversely, any informal proof provided in the core of the paper should be comple-
879 mented by formal proofs provided in appendix or supplemental material.
- 880 • Theorems and Lemmas that the proof relies upon should be properly referenced.

881 **4. Experimental result reproducibility**

882 Question: Does the paper fully disclose all the information needed to reproduce the main
883 experimental results of the paper to the extent that it affects the main claims and/or conclu-
884 sions of the paper (regardless of whether the code and data are provided or not)?

885 Answer: [Yes]

886 Justification: The paper fully discloses all the necessary information to reproduce the main
887 experimental results, including detailed descriptions of the methodologies, parameters, and
888 conditions that affect the main claims and conclusions.

889 Guidelines:

- 890 • The answer NA means that the paper does not include experiments.
- 891 • If the paper includes experiments, a No answer to this question will not be perceived
892 well by the reviewers: Making the paper reproducible is important, regardless of
893 whether the code and data are provided or not.
- 894 • If the contribution is a dataset and/or model, the authors should describe the steps
895 taken to make their results reproducible or verifiable.
- 896 • Depending on the contribution, reproducibility can be accomplished in various ways.
897 For example, if the contribution is a novel architecture, describing the architecture
898 fully might suffice, or if the contribution is a specific model and empirical evaluation,
899 it may be necessary to either make it possible for others to replicate the model with
900 the same dataset, or provide access to the model. In general, releasing code and data
901 is often one good way to accomplish this, but reproducibility can also be provided via
902 detailed instructions for how to replicate the results, access to a hosted model (e.g., in
903 the case of a large language model), releasing of a model checkpoint, or other means
904 that are appropriate to the research performed.
- 905 • While NeurIPS does not require releasing code, the conference does require all sub-
906 missions to provide some reasonable avenue for reproducibility, which may depend
907 on the nature of the contribution. For example
 - 908 (a) If the contribution is primarily a new algorithm, the paper should make it clear
909 how to reproduce that algorithm.
 - 910 (b) If the contribution is primarily a new model architecture, the paper should describe
911 the architecture clearly and fully.
 - 912 (c) If the contribution is a new model (e.g., a large language model), then there should
913 either be a way to access this model for reproducing the results or a way to re-
914 produce the model (e.g., with an open-source dataset or instructions for how to
915 construct the dataset).
 - 916 (d) We recognize that reproducibility may be tricky in some cases, in which case au-
917 thors are welcome to describe the particular way they provide for reproducibility.
918 In the case of closed-source models, it may be that access to the model is limited in
919 some way (e.g., to registered users), but it should be possible for other researchers
920 to have some path to reproducing or verifying the results.

921 **5. Open access to data and code**

922 Question: Does the paper provide open access to the data and code, with sufficient instruc-
923 tions to faithfully reproduce the main experimental results, as described in supplemental
924 material?

925 Answer: [No]

926 Justification: The paper does not currently provide open access to the data and code, but
927 this will be addressed and included later.

928 Guidelines:

- 929 • The answer NA means that paper does not include experiments requiring code.
- 930 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 931 • While we encourage the release of code and data, we understand that this might not
932 be possible, so No is an acceptable answer. Papers cannot be rejected simply for not
933 including code, unless this is central to the contribution (e.g., for a new open-source
934 benchmark).
- 935 • The instructions should contain the exact command and environment needed to run to
936 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 937 • The authors should provide instructions on data access and preparation, including how
938 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 939 • The authors should provide scripts to reproduce all experimental results for the new
940 proposed method and baselines. If only a subset of experiments are reproducible, they
941 should state which ones are omitted from the script and why.
- 942 • At submission time, to preserve anonymity, the authors should release anonymized
943 versions (if applicable).
- 944 • Providing as much information as possible in supplemental material (appended to the
945 paper) is recommended, but including URLs to data and code is permitted.

946 6. Experimental setting/details

947 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
948 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
949 results?

950 Answer: [Yes]

951 Justification: In the description of the experiment, this article specifies all the training and
952 test details.

953 Guidelines:

- 954 • The answer NA means that the paper does not include experiments.
- 955 • The experimental setting should be presented in the core of the paper to a level of
956 detail that is necessary to appreciate the results and make sense of them.
- 957 • The full details can be provided either with the code, in appendix, or as supplemental
958 material.

959 7. Experiment statistical significance

960 Question: Does the paper report error bars suitably and correctly defined or other appropri-
961 ate information about the statistical significance of the experiments?

962 Answer: [No]

963 Justification: The experiments in this paper are not suitable for quantitative statistical anal-
964 ysis, so there are no error bars or similar components.

965 Guidelines:

- 966 • The answer NA means that the paper does not include experiments.
- 967 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
968 dence intervals, or statistical significance tests, at least for the experiments that support
969 the main claims of the paper.

- 972 • The factors of variability that the error bars are capturing should be clearly stated (for
 973 example, train/test split, initialization, random drawing of some parameter, or overall
 974 run with given experimental conditions).
 975 • The method for calculating the error bars should be explained (closed form formula,
 976 call to a library function, bootstrap, etc.)
 977 • The assumptions made should be given (e.g., Normally distributed errors).
 978 • It should be clear whether the error bar is the standard deviation or the standard error
 979 of the mean.
 980 • It is OK to report 1-sigma error bars, but one should state it. The authors should prefer-
 981 ably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of
 982 Normality of errors is not verified.
 983 • For asymmetric distributions, the authors should be careful not to show in tables or
 984 figures symmetric error bars that would yield results that are out of range (e.g. negative
 985 error rates).
 986 • If error bars are reported in tables or plots, The authors should explain in the text how
 987 they were calculated and reference the corresponding figures or tables in the text.

988 **8. Experiments compute resources**

989 Question: For each experiment, does the paper provide sufficient information on the com-
 990 puter resources (type of compute workers, memory, time of execution) needed to reproduce
 991 the experiments?

992 Answer: [Yes]

993 Justification: The paper provides sufficient information on the computer resources required
 994 to reproduce each experiment

995 Guidelines:

- 996 • The answer NA means that the paper does not include experiments.
 997 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
 998 or cloud provider, including relevant memory and storage.
 999 • The paper should provide the amount of compute required for each of the individual
 1000 experimental runs as well as estimate the total compute.
 1001 • The paper should disclose whether the full research project required more compute
 1002 than the experiments reported in the paper (e.g., preliminary or failed experiments
 1003 that didn't make it into the paper).

1004 **9. Code of ethics**

1005 Question: Does the research conducted in the paper conform, in every respect, with the
 1006 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

1007 Answer: [Yes]

1008 Justification: The research conducted in the paper conforms to the NeurIPS Code of Ethics
 1009 in all respects, ensuring that ethical considerations were prioritized throughout the study.

1010 Guidelines:

- 1011 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
 1012 • If the authors answer No, they should explain the special circumstances that require a
 1013 deviation from the Code of Ethics.
 1014 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
 1015 eration due to laws or regulations in their jurisdiction).

1016 **10. Broader impacts**

1017 Question: Does the paper discuss both potential positive societal impacts and negative
 1018 societal impacts of the work performed?

1019 Answer: [Yes]

1020 Justification: The discussion is included in the sections on Conclusion and Limitations of
 1021 the paper.

1022 Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The paper poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: The creators or original owners of the assets used in the paper are properly credited, and the license and terms of use are explicitly mentioned and respected.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.

- 1075 • For scraped data from a particular source (e.g., website), the copyright and terms of
1076 service of that source should be provided.
1077 • If assets are released, the license, copyright information, and terms of use in the pack-
1078 age should be provided. For popular datasets, paperswithcode.com/datasets has
1079 curated licenses for some datasets. Their licensing guide can help determine the li-
1080 cense of a dataset.
1081 • For existing datasets that are re-packaged, both the original license and the license of
1082 the derived asset (if it has changed) should be provided.
1083 • If this information is not available online, the authors are encouraged to reach out to
1084 the asset's creators.

1085 **13. New assets**

1086 Question: Are new assets introduced in the paper well documented and is the documenta-
1087 tion provided alongside the assets?

1088 Answer: [No]

1089 Justification: The paper does not currently provide open access to the new assets, but this
1090 will be addressed and included later.

1091 Guidelines:

- 1092 • The answer NA means that the paper does not release new assets.
1093 • Researchers should communicate the details of the dataset/code/model as part of their
1094 submissions via structured templates. This includes details about training, license,
1095 limitations, etc.
1096 • The paper should discuss whether and how consent was obtained from people whose
1097 asset is used.
1098 • At submission time, remember to anonymize your assets (if applicable). You can
1099 either create an anonymized URL or include an anonymized zip file.

1100 **14. Crowdsourcing and research with human subjects**

1101 Question: For crowdsourcing experiments and research with human subjects, does the pa-
1102 per include the full text of instructions given to participants and screenshots, if applicable,
1103 as well as details about compensation (if any)?

1104 Answer: [NA]

1105 Justification: The paper does not involve crowdsourcing nor research with human subjects.

1106 Guidelines:

- 1107 • The answer NA means that the paper does not involve crowdsourcing nor research
1108 with human subjects.
1109 • Including this information in the supplemental material is fine, but if the main contri-
1110 bution of the paper involves human subjects, then as much detail as possible should
1111 be included in the main paper.
1112 • According to the NeurIPS Code of Ethics, workers involved in data collection, cura-
1113 tion, or other labor should be paid at least the minimum wage in the country of the
1114 data collector.

1115 **15. Institutional review board (IRB) approvals or equivalent for research with human
1116 subjects**

1117 Question: Does the paper describe potential risks incurred by study participants, whether
1118 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
1119 approvals (or an equivalent approval/review based on the requirements of your country or
1120 institution) were obtained?

1121 Answer: [NA]

1122 Justification: The paper does not involve crowdsourcing nor research with human subjects.

1123 Guidelines:

- 1124 • The answer NA means that the paper does not involve crowdsourcing nor research
1125 with human subjects.

- 1126 • Depending on the country in which research is conducted, IRB approval (or equivalent)
1127 may be required for any human subjects research. If you obtained IRB approval,
1128 you should clearly state this in the paper.
1129 • We recognize that the procedures for this may vary significantly between institutions
1130 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
1131 guidelines for their institution.
1132 • For initial submissions, do not include any information that would break anonymity
1133 (if applicable), such as the institution conducting the review.

1134 **16. Declaration of LLM usage**

1135 Question: Does the paper describe the usage of LLMs if it is an important, original, or
1136 non-standard component of the core methods in this research? Note that if the LLM is used
1137 only for writing, editing, or formatting purposes and does not impact the core methodology,
1138 scientific rigorousness, or originality of the research, declaration is not required.

1139 Answer: [Yes]

1140 Justification: The paper describes the usage of LLMs as a crucial component of the core
1141 methods, highlighting their importance and originality in the research.

1142 Guidelines:

- 1143 • The answer NA means that the core method development in this research does not
1144 involve LLMs as any important, original, or non-standard components.
1145 • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>)
1146 for what should or should not be described.