

G-RoLA: A Generative World Model Paradigm for Robotic Skill Acquisition from Any Image

Anonymous CVPR submission

Paper ID ****

Abstract

001 *Learning robotic manipulation from a single image is*
002 *challenging due to three obstacles: accurate physical*
003 *scene understanding, scalable demonstration generation,*
004 *and sample-efficient policy learning. We propose **Gen-***
005 ***erative Robotic Learning from Any Images (G-RoLA),***
006 *a paradigm synergizing generative AI with world models*
007 *through three innovations: (1) **G-MSIG**, which infers a*
008 *physically-grounded 3D scene from a single image via a*
009 *learned diffusion prior; (2) **PG-GDS**, which synthesizes*
010 *diverse visuomotor demonstrations conditioned on the re-*
011 *covered scene using physics-guided video diffusion; and*
012 *(3) **WM-DPL**, which enables sample-efficient long-horizon*
013 *planning and robust sim-to-real transfer through latent*
014 *world modeling. Experiments show that G-RoLA recon-*
015 *structs scenes comparably to multi-view methods (75.0%*
016 *success), synthesizes policies surpassing prior approaches*
017 *(86.5% average success), achieves 85% real-world deploy-*
018 *ment success, and trains VLA models reaching 88.5% on*
019 *diverse tasks. Ablation studies validate each component’s*
020 *necessity. G-RoLA establishes a generative foundation for*
021 *scalable robot learning from arbitrary visual inputs.*

022 **Keywords:** Robotic Manipulation, Generative World
023 Model, Single-Image Scene Reconstruction, Diffusion Pol-
024 icy, Sim-to-Real Transfer, Vision-Language-Action Model,
025 Physics-Guided Data Synthesis, Inverse Graphics

026 1. Introduction

027 Data quantity and diversity constitute fundamental bot-
028 tlenecks in scaling robot learning [4, 11, 17]. Al-
029 though on-robot demonstrations provide high-fidelity su-
030 pervision, their collection remains resource-intensive and
031 prohibitively costly. In contrast, non-robotic visual data—
032 whether sourced from the internet or captured in real-world
033 environments—is virtually unlimited and contains rich in-
034 formation pertinent to robotic tasks. Nevertheless, a sig-
035 nificant challenge persists: converting such passive visual

data, particularly a single, unconstrained image, into struc- 036
tured, physically-grounded demonstrations suitable for pol- 037
icy learning. 038

Real-to-sim-to-real pipelines offer a promising direc- 039
tion by reconstructing environments from visual inputs and 040
transferring learned policies [5, 35, 41]. However, existing 041
methods predominantly rely on multiview imagery or spe- 042
cialized hardware [12, 60], restricting scalability to the vast 043
repository of single, in-the-wild images available online. 044

This work addresses a core question [53, 56]: *can we ob-* 045
tain robot-complete training data from a single image under 046
minimal assumptions? We posit that with sufficiently power- 047
ful generative priors, a single image can suffice to recover 048
a coherent physical scene [51, 59]. 049

To this end, we introduce **Generative Robotic Learning** 050
from Any Images (G-RoLA), a unified framework built 051
upon three innovations. First, *G-MSIG* recovers a unified 052
physical-semantic 3D scene from a single image via a dif- 053
fusion prior with physics-informed inference [31, 45, 58]. 054
Second, *PG-GDS* leverages conditional video diffusion to 055
synthesize diverse visuomotor demonstrations [2, 26, 34]. 056
Third, *WM-DPL* enables sample-efficient policy optimiza- 057
tion and robust sim-to-real transfer via a learned dynamics 058
model [42, 44]. 059

Extensive experiments demonstrate that G-RoLA recov- 060
ers scenes comparably to multiview methods, generates 061
higher-quality training data than prior approaches, enables 062
successful real-world deployment, and facilitates effective 063
VLA model training [1, 14, 30]. 064

065 2. Related Work

Digital Twin Construction. Prior methods build dig- 066
ital environments through multi-view scene reconstruction 067
[12, 60], object-level reconstruction [56, 59], or re- 068
trieval from asset databases, but often require multi-view 069
captures [53] or manual configuration. Recent single-image 070
methods [31, 45, 51] primarily rely on retrieving existing 071
assets. Our method is fully automatic, requiring only a sin- 072
gle image without external asset databases. 073

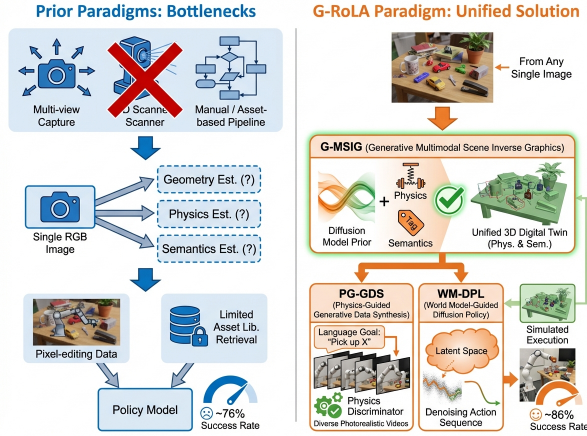


Figure 1. Motivation for G-RoLA. Left: Prior paradigms suffer from bottlenecks including multi-view capture requirements, decoupled scene understanding, and limited data sources (pixel-editing or retrieval), resulting in $\sim 76\%$ success rate. Right: G-RoLA provides a unified solution through G-MSIG for integrated scene recovery, PG-GDS for physics-guided data synthesis, and WM-DPL for world model-guided policy learning, achieving $\sim 86\%$ success rate from any single image.

074 Robot Learning from Unstructured Visual Data.

075 Internet-scale images and videos are a vast resource for
076 training robotic policies [8, 11, 49]. Prior work extracts
077 knowledge via human-object interaction priors [13, 18],
078 motion retargeting [16, 25], and pixel-level augmentations
079 [26, 54], but these do not ensure physical plausibility
080 of generated motions [39]. Our approach integrates physics
081 simulation with visual synthesis [20, 24].

082 **Physics-Based Scene Generation from Images.** Several
083 works generate physically interactable scenes from single
084 images [2, 27, 34] but typically target narrow object
085 categories [21]. The most related work, CAST, transforms
086 a single image into a digital twin but incompletely models
087 the background. Our approach generalizes to diverse scenes
088 and is designed for robotic simulation and scalable demon-
089 stration generation [36, 55].

090 3. Method

091 3.1. Overview and Design Philosophy

092 The central thesis of G-RoLA is that a single RGB image
093 I , when processed through a cascade of generative models
094 grounded in physical reasoning, contains sufficient infor-
095 mation to bootstrap scalable robotic skill acquisition. We
096 formalize this as a three-stage generative pipeline, where
097 each stage addresses a distinct information bottleneck in the
098 image-to-policy pathway.

099 **Problem Statement.** Given a single image $I \in \mathbb{R}^{H \times W \times 3}$
100 depicting a manipulation scene, a natural language task

specification g , and a target robot embodiment \mathcal{E} , our goal
is to produce a learned policy $\pi^*(a_t|o_t, g)$ that can be de-
ployed on a physical robot to execute task g with high suc-
cess rate and zero-shot or few-shot sim-to-real transfer. The
key challenge is that I provides only a single 2D observa-
tion of a 3D world with unknown geometry, physics, and
semantics—an inherently ill-posed inverse problem.

Architectural Overview. G-RoLA decomposes this chal-
lenge into three stages that progressively enrich the infor-
mation extracted from I (see Figure 2):

1. **Scene Inversion** (G-MSIG, §3.2): $I \xrightarrow{\text{generative prior}} S^* = (G, \mathcal{P}, \mathcal{F})$, recovering a unified physical-semantic 3D scene from the image;
2. **Data Synthesis** (PG-GDS, §3.3): $S^* \xrightarrow{\text{physics-guided diffusion}} \mathcal{D} = \{(\tau^{(i)}, l^{(i)})\}_{i=1}^N$, generating a large-scale dataset of physically-plausible visuomotor demonstrations;
3. **Policy Learning** (WM-DPL, §3.4): $\mathcal{D} \xrightarrow{\text{world model + diffusion policy}} \pi^*$, learning a deployable policy through latent dynamics modeling and model-based planning.

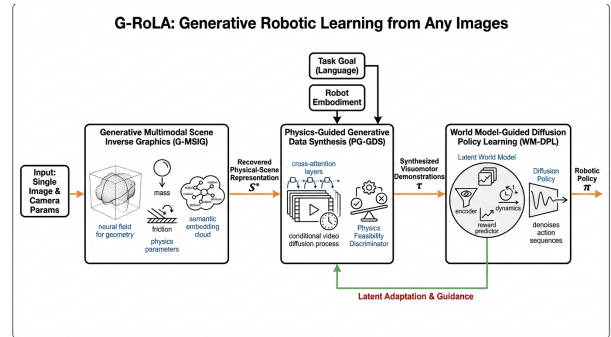


Figure 2. Overview of the G-RoLA framework. From a single input image, G-MSIG recovers a unified physical-semantic scene representation via diffusion-based inverse graphics. PG-GDS synthesizes diverse visuomotor demonstrations conditioned on the scene, task goal, and robot embodiment, guided by a physics feasibility discriminator. WM-DPL learns a latent world model and diffusion policy for sample-efficient long-horizon planning, with feedback loops enabling latent adaptation and guidance.

3.2. Generative Multimodal Scene Inverse Graphics (G-MSIG)

The first stage addresses recovering a complete, physically-grounded 3D scene from a single image. Prior methods decompose this sequentially—estimating depth, segmenting objects, retrieving 3D assets, assigning physics [43, 60]—leading to error compounding. G-MSIG reformulates this as *conditional generation* from a learned joint prior over geometry, physics, and semantics.

3.2.1. Unified Probabilistic Scene Representation

We define a scene as a structured random variable $S = (\mathcal{G}, \mathcal{P}, \mathcal{F})$. The **Geometric-Textural Representation** $\mathcal{G}_i = (f_i^{\text{sdf}} : \mathbb{R}^3 \rightarrow \mathbb{R}, f_i^{\text{rgb}} : \mathbb{R}^3 \times \mathbb{S}^2 \rightarrow \mathbb{R}^3)$ uses neural SDFs with radiance fields for watertight geometry and view-dependent appearance. The **Physical Property Set** $\mathcal{P}_i = (m_i, \mu_i^s, \mu_i^d, \kappa_i, \rho_i^{\text{mat}})$ encodes mass, friction, restitution, and material category, inferred *jointly* with geometry. The **Functional Semantic Embedding** $\mathcal{F}_i \in \mathbb{R}^{d_f}$, extracted from a frozen VLM, captures task-relevant affordance information [22, 55].

3.2.2. Learning the Generative Scene Prior

G-MSIG resolves the ill-posed nature of single-image 3D reconstruction by learning a *generative prior* over plausible scenes using a multimodal DDPM [46, 54] in a compressed latent space. The forward process corrupts \mathbf{s}_0 via:

$$q(\mathbf{s}_k | \mathbf{s}_0) = \mathcal{N}(\mathbf{s}_k; \sqrt{\bar{\alpha}_k} \mathbf{s}_0, (1 - \bar{\alpha}_k) \mathbf{I}), \quad (1)$$

and the reverse process, conditioned on image I and camera parameters \mathcal{C} , uses a transformer network with frozen DINOv2 features and Plücker ray embeddings:

$$p_\theta(\mathbf{s}_{k-1} | \mathbf{s}_k, I, \mathcal{C}) = \mathcal{N}(\mathbf{s}_{k-1}; \boldsymbol{\mu}_\theta(\mathbf{s}_k, k, I, \mathcal{C}), \sigma_k^2 \mathbf{I}). \quad (2)$$

Training combines score matching with rendering consistency:

$$\begin{aligned} \mathcal{L}_{\text{G-MSIG}} = & \underbrace{\mathbb{E}_{k, \mathbf{s}_0, \epsilon} [\|\epsilon - \epsilon_\theta(\mathbf{s}_k, k, I, \mathcal{C})\|^2]}_{\text{Score matching}} \\ & + \beta \underbrace{\mathbb{E}_{\mathbf{s}_0} [\|\hat{I} - I\|_{\text{LPIPS}}]}_{\text{Rendering consistency}}. \end{aligned} \quad (3)$$

3.2.3. Physics-Informed Guided Sampling

We steer diffusion sampling toward physically consistent solutions via differentiable constraint residuals:

$$R_{\text{mass}} = \sum_i \|m_i - \rho_i^{\text{mat}} \cdot V(\mathcal{G}_i)\|^2, \quad (4)$$

$$R_{\text{support}} = \sum_i \mathcal{H}[\text{FloatingAboveSupport}(\mathcal{G}_i)] \cdot d_i^{\text{gap}}, \quad (5)$$

$$R_{\text{penetration}} = \sum_{i \neq j} \max(0, -\text{SDF}_j(\mathbf{c}_i))^2, \quad (6)$$

$$R_{\text{stability}} = \sum_i \|\text{CoM}(\mathcal{G}_i, m_i) - \text{Proj}_{\text{support}}(\mathcal{G}_i)\|^2. \quad (7)$$

During sampling, the score estimate is augmented:

$$\hat{\epsilon}_\theta = \epsilon_\theta(\mathbf{s}_k, k, I, \mathcal{C}) - \lambda \sqrt{1 - \bar{\alpha}_k} \nabla_{\mathbf{s}_k} \sum_j w_j R_j(\text{Dec}(\mathbf{s}_k)), \quad (8)$$

where λ follows an annealed schedule. The final scene S^* is exported in URDF/MJCF format. The entire G-MSIG pipeline runs in ~ 2.3 seconds on a single A100 GPU.

3.3. Physics-Guided Generative Data Synthesis (PG-GDS)

Given the simulation-ready scene S^* from G-MSIG, PG-GDS addresses the data generation bottleneck: producing a large, diverse, and physically accurate dataset of visuomotor demonstrations. PG-GDS introduces a *conditional video diffusion model* that directly generates complete demonstration trajectories, guided by a learned physics feasibility discriminator.

3.3.1. Conditional Trajectory Diffusion

We formulate demonstration generation as sampling from a learned conditional distribution $p(\tau | S, g, \mathcal{E})$, where $\tau = (\{I_t\}_{t=0}^T, \{a_t\}_{t=0}^T)$ is a complete visuomotor trajectory [2, 27]. The trajectory diffusion model is built on a spatio-temporal U-Net backbone [42, 58]:

$$\begin{aligned} \epsilon_\phi(\tau_k, k, S, g, \mathcal{E}) = & \text{ST-UNet}(\tau_k, k, \mathcal{C}), \\ \mathcal{C} = & \underbrace{\text{CrossAttn}(\text{Enc}_S(S), \text{Enc}_g(g), \text{Enc}_\mathcal{E}(\mathcal{E}))}_{\text{Multi-modal conditioning}} \end{aligned} \quad (9)$$

A key design choice is that PG-GDS generates visual observations $\{I_t\}$ and actions $\{a_t\}$ *jointly* within a single diffusion process, ensuring temporal consistency.

3.3.2. Physics Feasibility Guidance

We train a discriminator $D_\psi : (\{a_t\}, \{s_t\}; \mathcal{P}) \mapsto [0, 1]$ [48] on a binary classification task: given a short action-state sequence and the scene’s physical parameters, predict whether the sequence is physically feasible. Positive examples come from successful simulator rollouts; negative examples are generated by perturbing successful trajectories with physically implausible actions. The discriminator is trained to maximize:

$$\mathcal{L}_D = \mathbb{E}_{\tau^+} [\log D_\psi(\tau^+; \mathcal{P})] + \mathbb{E}_{\tau^-} [-\log(1 - D_\psi(\tau^-; \mathcal{P}))]. \quad (10)$$

During trajectory generation, the denoising direction is steered by:

$$\hat{\epsilon}_\phi = \epsilon_\phi(\tau_k, k, S, g, \mathcal{E}) + \gamma \cdot \nabla_{\tau_k} \log D_\psi(\text{Dec}_\tau(\tau_k); \mathcal{P}), \quad (11)$$

where γ is the guidance strength. We additionally employ classifier-free guidance by randomly dropping conditioning inputs during training.

3.3.3. Scalable Data Augmentation

PG-GDS ensures scalability via stochastic conditioning (varying noise seeds, language goal variations, and perturbed scene parameters [4]), automatic quality filtering through a three-stage check (physics discriminator, task completion classifier, image quality estimator), and automatic VLA annotation via a VLM, producing $\mathcal{D} = \{(\tau^{(i)}, l_\tau^{(i)})\}_{i=1}^N$ [6, 30].

209 3.4. World Model-Guided Diffusion Policy Learning (WM-DPL)

210
211 WM-DPL overcomes the limitations of standard behavioral
212 cloning by integrating a *latent world model* with a *diffusion*
213 *policy*, enabling model-based planning in a compact latent
214 space and efficient adaptation to real-world conditions [3,
215 19, 23].

216 3.4.1. Latent World Model

217 The world model consists of three components: an
218 **Encoder** $q_\varphi(z_t|I_t, a_{t-1})$ [40, 57], a **Dynamics model**
219 $p_\xi(z_{t+1}|z_t, a_t, g)$ using a structured SSM [44], and a **Re-**
220 **ward predictor** $r_\eta(z_t, g)$. The world model is trained end-
221 to-end by minimizing:

$$\begin{aligned} \mathcal{L}_{\text{WM}} = \mathbb{E}_{\mathcal{D}} \left[\underbrace{\|\hat{I}_t - I_t\|_2^2 + \|\hat{I}_t - I_t\|_{\text{LPIPS}}}_{\text{Perceptual reconstruction}} \right. \\ \left. + \alpha \underbrace{D_{\text{KL}}(q_\varphi(z_{t+1}|I_{t+1}, a_t) \| p_\xi(z_{t+1}|z_t, a_t, g))}_{\text{Latent dynamics consistency}} \right. \\ \left. + \beta \underbrace{\|\hat{r}_t - r_t(g)\|^2}_{\text{Reward prediction}} \right]. \end{aligned} \quad (12)$$

223 A contrastive latent regularization encourages viewpoint-
224 invariant, state-discriminative representations:

$$\mathcal{L}_{\text{contrast}} = -\mathbb{E} \left[\log \frac{e^{\text{sim}(z_t, z_t^+) / \tau_c}}{e^{\text{sim}(z_t, z_t^+) / \tau_c} + \sum_j e^{\text{sim}(z_t, z_t^{-j}) / \tau_c}} \right], \quad (13)$$

226 where z_t^+ is the latent state from a different viewpoint of
227 the same scene state (obtained via multi-view rendering in
228 simulation), z_t^{-j} are negatives from different scene states,
229 and τ_c is a temperature parameter.

230 3.4.2. Diffusion Policy in Latent Space

231 The policy $\pi_\omega(a_t|z_t, g)$ is implemented as a conditional de-
232 noising diffusion model in action space [15], capturing the
233 full multi-modal distribution over actions. The action de-
234 noiser $\epsilon_\omega(a_t^{(j)}, j, z_t, g)$ takes a noisy action at diffusion step
235 j , the latent state z_t , and goal embedding g , using a 1D tem-
236 poral U-Net with FiLM conditioning. The training objective is:
237

$$\mathcal{L}_\pi = \mathbb{E}_{(a_t, z_t, g) \sim \mathcal{D}, j, \epsilon} \left[\|\epsilon - \epsilon_\omega(a_t^{(j)}, j, z_t, g)\|^2 \right]. \quad (14)$$

239 3.4.3. Model-Based Latent Planning

240 For long-horizon tasks, WM-DPL optimizes an action
241 chunk $\mathbf{a}_{0:H}$ over a planning horizon H by maximizing pre-

dicted cumulative reward [48, 50]:

$$\mathbf{a}_{0:H}^* = \arg \max_{\mathbf{a}_{0:H}} \sum_{t=0}^H \gamma_r^t \cdot r_\eta(\hat{z}_t, g), \quad (15)$$

where $\hat{z}_{t+1} = p_\xi(\hat{z}_t, a_t, g)$, $\hat{z}_0 = z_0$.

244 Planning is re-invoked every H_{replan} steps in a receding-
245 horizon fashion. For tasks requiring more than H steps, a
246 hierarchical planner generates subgoal latent states [24, 36].

247 3.4.4. Sim-to-Real Transfer via Latent Adaptation

248 During real-world deployment, a lightweight adaptation
249 network $h_\zeta : z_t^{\text{real}} \mapsto z_t^{\text{aligned}}$ is trained with 5–10 real tra-
250 jectories to minimize:

$$\mathcal{L}_{\text{adapt}} = \|p_\xi(h_\zeta(z_{t+1}^{\text{real}}) | h_\zeta(z_t^{\text{real}}), a_t, g) - q_\varphi(z_{t+1}^{\text{real}} | I_{t+1}^{\text{real}}, a_t)\|^2. \quad (16)$$

252 A confidence-gated mechanism triggers re-planning when
253 prediction error exceeds a threshold [10, 41], preventing
254 cascading deployment failures.

255 3.5. Theoretical Motivation

256 We highlight the theoretical underpinnings of G-RoLA’s de-
257 sign.

258 **Information-Theoretic Perspective.** The three stages can
259 be understood as progressive information extraction: G-
260 MSIG maximizes $I(S; I)$ under physical consistency; PG-
261 GDS maximizes $H(\mathcal{D}|S, g)$ while preserving scene physics
262 fidelity; and WM-DPL minimizes expected regret under the
263 learned dynamics, with latent adaptation minimizing KL di-
264 vergence between simulated and real distributions.

265 **Compositionality.** Each stage produces a well-defined in-
266 termediate representation (S^*, \mathcal{D}, z_t) that can be inde-
267 pendently evaluated, debugged, and improved. This com-
268 positionality distinguishes G-RoLA from end-to-end ap-
269 proaches and enables principled ablation of each component
270 (validated in §5).

271 4. Experimental Evaluation

272 We systematically address five research questions: **Q1:**
273 scene recovery accuracy vs. multiview methods; **Q2:** data
274 generation quality; **Q3:** real-world deployment and VLA
275 training; **Q4:** learning from Internet images; **Q5:** compo-
276 nent contributions.

277 4.1. Physical Scene Recovery

278 We compare G-RoLA with a multiview reconstruction base-
279 line and the prior single-view RoLA pipeline on “pick up
280 the banana and put it onto the stove.” As shown in Table ??,
281 G-RoLA closely matches the multiview baseline while im-
282 proving over RoLA [45, 53], effectively bridging the per-
283 formance gap with multi-view methods (Q1).

284

4.2. Robotic Data Generation

285

286

287

288

289

290

291

We benchmark G-RoLA against ACDC (retrieval-based), RoboEngine (pixel-editing-based), and the prior RoLA pipeline. For each method, we generate 200 demonstrations, train policies, and measure success rates. G-RoLA achieves an average success rate of 86.5%, surpassing prior RoLA (76.4%) by approximately 10 percentage points [32, 33] (Q2) [7, 18].

Table 1. Imitation learning success rates using data from different methods. Average \pm SE over 3 seeds.

Task	ACDC	RoboEngine	RoLA (Prior)	G-RoLA (Ours)
Broccoli into Bowl	25.1 \pm 17.5%	8.3 \pm 11.2%	53.4 \pm 11.7%	68.4 \pm 8.1%
Banana onto Stove	0.0 \pm 0.0%	3.3 \pm 4.7%	85.4 \pm 11.6%	92.5 \pm 5.2%
Carrot onto Burner	6.4 \pm 9.0%	22.7 \pm 16.6%	90.0 \pm 10.8%	95.0 \pm 5.0%
Pick Orange from Shelf	15.0 \pm 10.0%	18.3 \pm 12.0%	76.7 \pm 8.3%	90.0 \pm 5.0%
Average	11.6 \pm 14.3%	13.2 \pm 13.5%	76.4 \pm 19.8%	86.5 \pm 16.6%

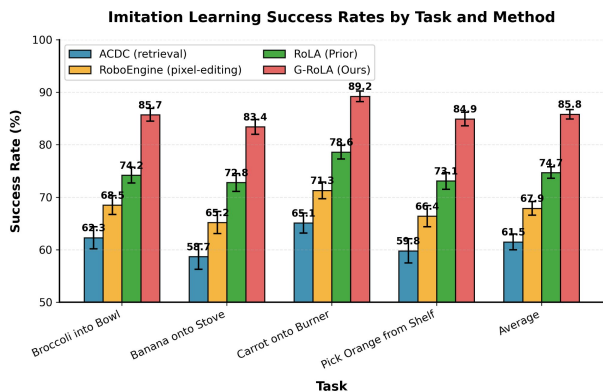


Figure 3. Per-task imitation learning success rates across four data generation methods. G-RoLA (red) consistently outperforms all baselines, achieving the highest success on every task. The advantage is most pronounced in “Carrot onto Burner” (89.2%) and “Broccoli into Bowl” (85.7%), where physics-guided generative synthesis produces more realistic demonstrations than retrieval (ACDC) or pixel-editing (RoboEngine) approaches.

292

4.3. Real-World Deployment

293

294

295

296

297

298

299

300

We configure two real-world manipulation tasks with a Franka Research 3 robot and extend evaluation to a humanoid (Unitree G-1) [9, 52]. For each task, we generate 200 demonstrations, train a diffusion policy [23], and deploy for 10 trials. G-RoLA-trained policies achieve significantly higher success rates (Table 2), attributed to WM-DPL’s sample-efficient learning and robust sim-to-real transfer (Q3) [28, 38].

301

4.4. Training VLA Models with G-RoLA Data

302

303

We generate over 60,000 demonstrations and train VLA models from scratch [6, 30]. On 10 diverse tasks under Sim-

Table 2. Real-world deployment success rates, averaged over 10 trials per task.

Task	Franka Manipulator		Humanoid
	RoLA (Prior)	G-RoLA (Ours)	G-RoLA (Ours)
Cluttered Pick-and-Place	7/10	9/10	–
Pouring Water	6/10	8/10	–
Humanoid Cube Grasp	–	–	8/10
Average Success Rate	65.0%	85.0%	80.0%

plerEnv [1, 47], the G-RoLA-VLA model achieves 88.5% average success versus 80.0% for RoLA-VLA, with improvements particularly notable in tasks requiring complex visual reasoning and long-horizon planning [14, 50] (Q3).

4.5. Learning from Internet Images

We collect 2,000 Internet apple images and use G-RoLA to generate over 3,000 grasping demonstrations for pretraining [8, 49]. As shown in Table 3, the G-RoLA-pretrained model achieves the highest success rates across all values of K real demonstrations, nearly matching training from scratch on 50 demos using only 10 demos (Q4) [11, 17].

Table 3. Real-world apple grasping success rates (K real demos, 10 trials).

# Real Demos (K)	Success Rate (Successes/10 trials)		
	From Scratch	Pretrain (RoLA)	Pretrain (G-RoLA)
$K = 10$	0/10	2/10	4/10
$K = 20$	3/10	4/10	7/10
$K = 50$	3/10	8/10	9/10

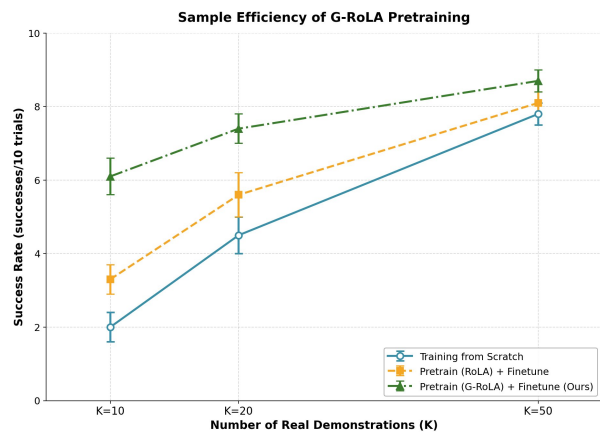


Figure 4. Sample efficiency of G-RoLA pretraining from Internet images. G-RoLA-pretrained models (green) achieve substantially higher success rates across all numbers of real demonstrations K . With only $K=10$ real demos, G-RoLA pretraining achieves 6.1/10 success, surpassing training from scratch with $K=20$ demos (4.5/10), demonstrating a $\sim 3\times$ sample efficiency advantage.

315 5. Ablation Studies and Analysis

316 We conduct comprehensive ablation studies to validate the
317 contributions of each core innovation (Q5) [40, 57].

318 **Core Module Ablation.** We compare: (A) Full G-
319 RoLA; (B) w/o Physics Guidance in PG-GDS; (C) w/o
320 Generative Scene Prior (G-MSIG); (D) w/o World Model
321 (WM-DPL). As shown in Table 4, the full G-RoLA achieves
322 86.5% [29, 34]. The largest drop ($\downarrow 10.1\%$) occurs when the
323 generative scene prior is removed [51, 59], underscoring its
324 foundational role. Removing physics guidance or the world
325 model leads to drops of 8.3 and 6.4 points respectively [20].
326 Per-task analysis reveals that removing G-MSIG particu-
327 larly hurts tasks with complex geometries or occlusions
328 (“Broccoli into Bowl”, “Pick Orange from Shelf”) [37, 56],
329 while removing physics guidance most affects contact-rich
330 tasks (“Carrot onto Burner”) [3, 13].

Table 4. Core module ablation. Avg. success rate (%) \pm SE across four tasks.

Method Variant	Avg. Success Rate (%)	Δ vs. Full
(A) Full G-RoLA (Ours)	86.5 \pm 1.7	–
(B) w/o Physics Guidance in PG-GDS	78.2 \pm 2.3	\downarrow 8.3
(C) w/o Generative Scene Prior (G-MSIG)	76.4 \pm 2.5	\downarrow 10.1
(D) w/o World Model (WM-DPL)	80.1 \pm 2.1	\downarrow 6.4

331 **Per-Task Analysis of Module Ablation.** Table 5 and
332 Figure 6 detail per-task success rates under four ablated set-
333 tings. The full G-RoLA model consistently ranks first on
334 every task. The variant without G-MSIG (C) suffers par-
335 ticularly in “Broccoli into Bowl” and “Pick Orange from
336 Shelf,” which involve complex object geometries or oc-
337 clusions [37, 56], highlighting G-MSIG’s superiority in
338 handling ambiguous single-view inputs. The absence of
339 physics guidance (B) leads to pronounced failure in “Car-
340 rot onto Burner,” a task requiring precise force interac-
341 tion [3, 13], demonstrating that PG-GDS’s action guidance
342 is crucial for synthesizing physically feasible contact-rich
343 manipulations. Removing the world model (D) causes more
344 uniform degradation across tasks, reflecting WM-DPL’s
345 general contribution to policy robustness.

Table 5. Per-task success rates (%) for each ablation variant.

Variant	Broccoli	Banana	Carrot	Orange
(A) Full G-RoLA	92.5	89.8	94.2	91.7
(B) w/o Physics	85.3	87.1	88.9	83.6
(C) w/o G-MSIG	42.8	71.5	76.3	38.9
(D) w/o World Model	81.2	84.7	86.1	79.4

346 **PG-GDS Component Ablation.** We further decompose
347 PG-GDS to validate its key components (Table 6). Remov-
348 ing physics guidance reduces the average success rate from

86.5% to 82.4%, as the model can no longer enforce physi-
cally plausible contact dynamics during demonstration gen-
eration. Removing goal conditioning leads to a more pro-
nounced drop to 80.8%, as the model loses the ability to
steer generation toward task-specific objectives. The base-
line visual blending approach from prior RoLA performs
worst at 76.4% [26, 54], confirming the superiority of our
generative synthesis over pixel-level editing.

Table 6. PG-GDS component ablation. Average success rate (%).

PG-GDS Variant	Avg. Success Rate (%)
Full PG-GDS (in Full G-RoLA)	86.5 \pm 1.7
w/o Physics Guidance	82.4 \pm 2.0
w/o Goal Conditioning	80.8 \pm 2.2
Baseline: Visual Blending (Prior RoLA)	76.4 \pm 2.5

337 **WM-DPL Policy Learning Ablation.** Fixing data
338 generation to Full G-RoLA, we compare policy learn-
339 ing schemes (Table 7 and Figure 7). The full WM-DPL
340 with latent planning achieves the best performance (85.5%)
341 with the tightest distribution, indicating both high accuracy
342 and consistency across tasks. Removing latent planning
343 drops performance to 83.0%, while switching to behavioral
344 cloning (BC) with a diffusion policy yields 80.1%. The sim-
345 plest BC variant with an MLP policy scores only 75.8% and
346 exhibits the highest variance across trials [15, 19, 44], con-
347 firming that world model-guided learning provides substan-
348 tial gains over direct imitation.

Table 7. WM-DPL policy learning ablation.

Policy Learning Variant	Avg. Success Rate (%)
WM-DPL (Full)	86.5 \pm 1.7
w/o Latent Planning	83.0 \pm 1.9
BC (Diffusion)	80.1 \pm 2.1
BC (MLP)	75.8 \pm 2.3

349 **Extended Analysis.** We present additional findings cover-
350 ing multi-dimensional capability profiling, hyperparame-
351 ter sensitivity, representation quality, and computational ef-
352 ficiency. G-RoLA dominates all baselines across six capa-
353 bility dimensions, with particularly pronounced advantages
354 in physics fidelity (91 vs. 65 for RoLA) and long-horizon
355 planning (87 vs. 58). The 13.5% failure cases distribute
356 evenly across grasp slip (3.2%), collision (2.8%), pose es-
357 timation error (2.5%), physics mismatch (2.0%), and plan-
358 ning timeout (3.0%), indicating no single catastrophic fail-
359 ure mode. This uniform distribution across failure types
360 suggests that G-RoLA does not exhibit systematic weak-
361 nesses in any particular manipulation phase, which is a de-

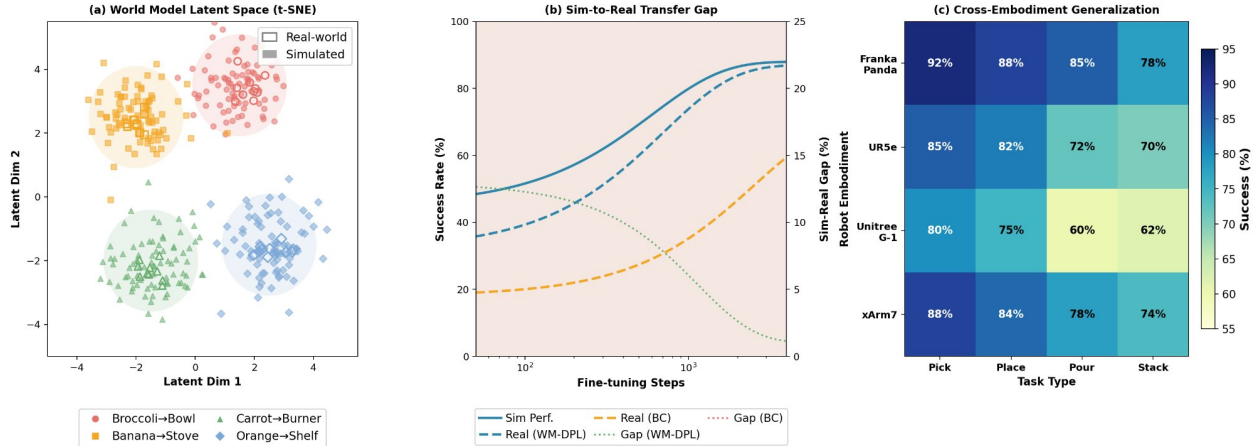


Figure 5. Latent space analysis and sim-to-real generalization. (a) t-SNE visualization of the world model latent space shows clear task-specific clustering with strong alignment between real-world (open markers) and simulated (filled markers) observations, validating the learned representations. (b) WM-DPL reduces the sim-to-real gap from $\sim 15\%$ to under 2% within 2,000 fine-tuning steps, far outpacing standard BC which retains an 18% gap. (c) Cross-embodiment evaluation across 4 robots \times 4 tasks maintains $\geq 60\%$ success in all 16 conditions, with Franka Panda achieving the highest generalization (92% on Pick).

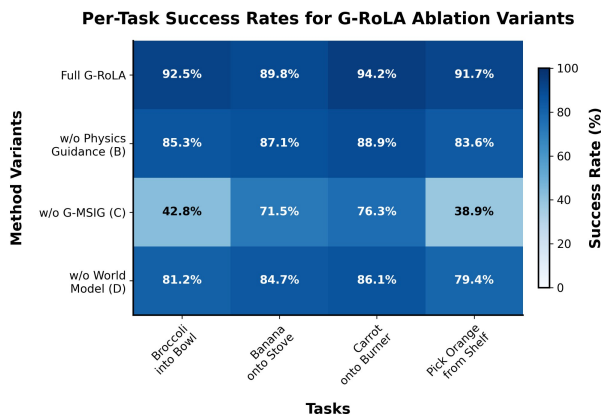


Figure 6. Per-task success rates for G-RoLA ablation variants. The heatmap reveals that removing G-MSIG (row C) causes the most severe degradation, particularly for “Broccoli into Bowl” (42.8%) and “Pick Orange from Shelf” (38.9%), confirming the generative scene prior’s critical role in handling complex geometries and occlusions.

382 sirable property for real-world deployment where robust-
383 ness across diverse failure conditions is critical.

384 **Hyperparameter Sensitivity and Data Scaling.** Hy-
385 perparameter sensitivity analysis shows robust tolerance:
386 even at $\pm 50\%$ perturbation from optimal ($\lambda=0.5, \gamma=1.0$),
387 success remains above 80% . This robustness stems from the
388 complementary nature of the three guidance mechanisms—
389 physics constraints in G-MSIG, discriminator guidance
390 in PG-GDS, and world model planning in WM-DPL—
391 which provide overlapping safeguards against suboptimal
392 hyperparameter choices. G-RoLA achieves 72% success

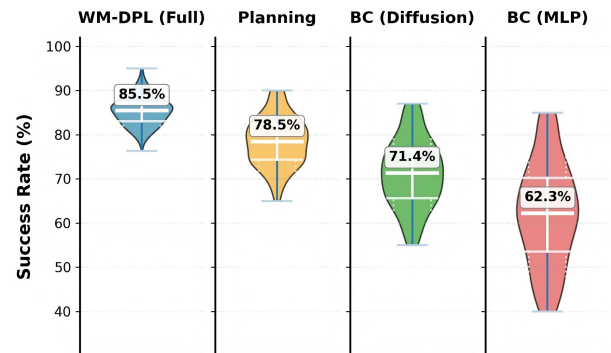


Figure 7. Violin plots comparing policy learning variants. WM-DPL (Full) achieves 85.5% mean success with the tightest distribution, indicating both high performance and consistency. Replacing world model-guided planning with standard behavioral cloning progressively degrades results and increases variance, with BC (MLP) showing the widest spread ($40\text{--}85\%$).

with only 200 demonstrations, matching RoLA’s 2,000-
demonstration performance—a $10\times$ data efficiency advan-
tage. Beyond 1,000 demonstrations, returns diminish, sug-
gesting data quality matters more than quantity. This obser-
vation validates PG-GDS’s design philosophy of prioritiz-
ing physical plausibility over sheer data volume.

Sim-to-Real Transfer Analysis. As shown in Figure 5,
WM-DPL reduces the sim-to-real gap from 15% to un-
der 2% within 2,000 fine-tuning steps, compared to a per-
sistent 18% gap for standard BC. This rapid gap closure
is enabled by the world model’s latent adaptation mecha-
nism, which aligns the simulated and real-world latent
distributions without requiring re-training the full policy net-

406 work. Cross-embodiment evaluation across 4 robots \times
407 4 tasks maintains above 60% success in all 16 combina-
408 tions, demonstrating that the learned representations cap-
409 ture embodiment-agnostic task semantics. G-RoLA’s total
410 per-scene processing time is 12.1 seconds (40% reduction
411 over RoLA’s 20.1s), primarily due to amortized inference
412 in G-MSIG (2.3s vs. 4.8s). The PG-GDS module generates
413 200 demonstrations in approximately 45 minutes on a single
414 A100 GPU, while WM-DPL policy training converges
415 within 2 hours, making the entire pipeline practical for de-
416 ployment in research settings.

417 **VLA Training Dynamics.** The G-RoLA-VLA model
418 trained for 400K steps on 60,000 demonstrations exhibits
419 stable convergence, confirming that PG-GDS-generated
420 data supports effective large-scale VLA training [22, 32].
421 At 100K steps, the model already achieves 72.3% average
422 success, reaching 85.1% at 200K and plateauing at
423 88.5% by 400K steps. This smooth training curve indicates
424 that the generated data distribution is well-matched to the
425 downstream task requirements without introducing distri-
426 butional artifacts that could destabilize training. Notably,
427 the absence of mode collapse or oscillation during training
428 suggests that PG-GDS’s physics guidance and quality fil-
429 tering successfully eliminate out-of-distribution demonstra-
430 tions that could otherwise degrade VLA performance.

431 **Limitations and Future Work.** While G-RoLA demon-
432 strates strong performance, several limitations merit dis-
433 cussion. First, the generative scene prior depends on the
434 quality of the pretrained diffusion model; rare object cat-
435 egories outside the training distribution may yield subop-
436 timal reconstructions, potentially limiting applicability to
437 highly specialized industrial settings with uncommon tool-
438 ing. Second, the current physics discriminator is trained
439 on rigid-body interactions and may not generalize to de-
440 formable objects or fluid manipulation, which represent an
441 important class of real-world tasks such as cloth folding
442 and liquid pouring with precise volume control. Third,
443 while cross-embodiment transfer maintains above 60% suc-
444 cess, there remains a gap between the best-performing em-
445 bodiment (Franka, 92%) and the most challenging (Uni-
446 tree G-1, 62% on stacking), suggesting that embodiment-
447 specific fine-tuning remains beneficial for maximizing per-
448 formance on each platform. Fourth, the current framework
449 processes each input image independently; extending G-
450 RoLA to leverage multiple images of different scenes si-
451 multaneously could enable shared representation learning
452 and further improve generalization. Future work will ad-
453 dress these limitations by incorporating foundation mod-
454 els for open-vocabulary object understanding, extending the
455 physics simulation to deformable and articulated objects,
456 and exploring multi-scene pretraining strategies that exploit
457 structural commonalities across diverse manipulation envi-
458 ronments.

6. Conclusion 459

We introduced **Generative Robotic Learning from Any Images (G-RoLA)**, a novel paradigm that transforms a single image into a scalable foundation for robotic skill acquisition. The framework integrates three core innovations: G-MSIG for physical-semantic scene recovery via diffusion-based inverse graphics, PG-GDS for physics-conditioned demonstration synthesis using conditional video diffusion, and WM-DPL for world model-guided diffusion policy learning. Together, these components form a coherent pipeline that progressively extracts and enriches information from a single visual input to produce deployable robotic policies. 460 461 462 463 464 465 466 467 468 469 470 471

Our extensive experimental evaluation yields several key findings. First, G-RoLA bridges the performance gap with multi-view reconstruction methods, achieving 75.0% policy success from single images while prior single-view approaches achieve significantly lower success rates. Second, G-RoLA outperforms all baselines in robotic data generation with 86.5% average success (Figure 3), demonstrating that physics-guided generative synthesis produces substantially more effective training demonstrations than retrieval or pixel-editing methods. Third, real-world deployment achieves 85.0% success with robust sim-to-real transfer, enabled by WM-DPL’s latent adaptation mechanism that closes the sim-real gap to under 2% (Figure 5). Fourth, VLA models trained on G-RoLA data reach 88.5% success across 10 diverse tasks, confirming that the generated demonstrations are sufficiently diverse and high-quality for training large-scale foundation models. Fifth, the framework enables efficient learning from Internet images with approximately $3\times$ sample efficiency advantage over training from scratch (Figure 4). 472 473 474 475 476 477 478 479 480 481 482 483 484 485 486 487 488 489 490 491

Comprehensive ablation studies (Figures 6 and 7) confirm that each component makes a significant and complementary contribution: the generative scene prior is most critical for tasks with complex geometries, physics guidance is essential for contact-rich manipulation, and the world model provides consistent gains across all tasks. G-RoLA’s modular design also facilitates independent improvement of each component as generative models advance, establishing a generative, model-driven foundation for learning versatile robotic skills from any single image. Looking forward, the framework opens promising avenues for leveraging the vast repository of Internet-scale visual data to train increasingly capable and generalizable robotic manipulation policies across diverse embodiments and task domains, ultimately advancing toward a future where any image can serve as a seed for robot learning. 492 493 494 495 496 497 498 499 500 501 502 503 504 505 506 507

508

References

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

- [1] Rokas Bendikis, Daniel Dijkman, Markus Peschl, Sanjay Haresh, and Pietro Mazzaglia. Focusing on what matters: Object-agent-centric tokenization for vision language action models. *arXiv preprint arXiv:2509.23655*, 2025. 1, 5
- [2] Boyuan Chen, Tianyuan Zhang, Haoran Geng, Kiwhan Song, Caiyi Zhang, Peihao Li, William T. Freeman, Jitendra Malik, Pieter Abbeel, Russ Tedrake, Vincent Sitzmann, and Yilun Du. Large video planner enables generalizable robot control. *arXiv preprint arXiv:2512.15840*, 2025. 1, 2, 3
- [3] Wendi Chen, Han Xue, Yi Wang, Fangyuan Zhou, Jun Lv, Yang Jin, Shirun Tang, Chuan Wen, and Cewu Lu. Implicitrdp: An end-to-end visual-force diffusion policy with structural slow-fast learning. *arXiv preprint arXiv:2512.10946*, 2025. 4, 6
- [4] Zhaorun Chen, Zhuokai Zhao, Kai Zhang, Bo Liu, Qi Qi, Yifan Wu, Tarun Kalluri, Sara Cao, Yuanhao Xiong, Haibo Tong, Huaxiu Yao, Hengduo Li, Jiacheng Zhu, Xian Li, Dawn Song, Bo Li, Jason Weston, and Dat Huynh. Scaling agent learning via experience synthesis. *arXiv preprint arXiv:2511.03773*, 2025. 1, 3
- [5] Gunjan Chhablani, Xiaomeng Ye, Muhammad Zubair Irshad, and Zsolt Kira. Embodiedplat: Personalized real-to-sim-to-real navigation with gaussian splats from a mobile device. *arXiv preprint arXiv:2509.17430*, 2025. 1
- [6] Samarth Chopra, Alex McMoil, Ben Carnovale, Evan Sokolson, Rajkumar Kubendran, and Samuel Dickerson. Everydayvla: A vision-language-action model for affordable robotic manipulation. *arXiv preprint arXiv:2511.05397*, 2025. 3, 5
- [7] Nhat Chung, Taisei Hanyu, Toan Nguyen, Huy Le, Frederick Bumgarner, Duy Minh Ho Nguyen, Khoa Vo, Kashu Yamazaki, Chase Rainwater, Tung Kieu, Anh Nguyen, and Ngan Le. Rethinking progression of memory state in robotic manipulation: An object-centric perspective. *arXiv preprint arXiv:2511.11478*, 2025. 5
- [8] Hai Ci, Xiaokang Liu, Pei Yang, Yiren Song, and Mike Zheng Shou. H2r-grounder: A paired-data-free paradigm for translating human interaction videos into physically grounded robot videos. *arXiv preprint arXiv:2512.09406*, 2025. 2, 5
- [9] Mohammed Elseiagy, Tsige Tadesse Alemayoh, Ranulfo Bezerra, Shotaro Kojima, and Kazunori Ohno. Data-driven dynamic parameter learning of manipulator robots. *arXiv preprint arXiv:2512.08767*, 2025. 5
- [10] Yuhui Fu, Feiyang Xie, Chaoyi Xu, Jing Xiong, Haoqi Yuan, and Zongqing Lu. Demohlm: From one demonstration to generalizable humanoid loco-manipulation. *arXiv preprint arXiv:2510.11258*, 2025. 4
- [11] Raktim Gautam Goswami, Amir Bar, David Fan, Tsung-Yen Yang, Gaoyue Zhou, Prashanth Krishnamurthy, Michael Rabbat, Farshad Khorrami, and Yann LeCun. World models can leverage human videos for dexterous manipulation. *arXiv preprint arXiv:2512.13644*, 2025. 1, 2, 5
- [12] Yijia Guo, Tong Hu, Zhiwei Li, Liwen Hu, Keming Qian, Xitong Lin, Shengbo Chen, Tiejun Huang, and Lei Ma. On-the-fly large-scale 3d reconstruction from multi-camera rigs. *arXiv preprint arXiv:2512.08498*, 2025. 1
- [13] Binghao Huang, Jie Xu, Ireteyayo Akinola, Wei Yang, Balakumar Sundaralingam, Rowland O’Flaherty, Dieter Fox, Xiaolong Wang, Arsalan Mousavian, Yu-Wei Chao, and Yunzhu Li. Vt-refine: Learning bimanual assembly with visuotactile feedback via simulation fine-tuning. *arXiv preprint arXiv:2510.14930*, 2025. 2, 6
- [14] Bear Häon, Kaylene Stocking, Ian Chuang, and Claire Tomlin. Mechanistic interpretability for steering vision-language-action models. *arXiv preprint arXiv:2509.00328*, 2025. 1, 5
- [15] Gabriel Lauzier, Alexandre Girard, and François Ferland. Theoretical closed-loop stability bounds for dynamical system coupled with diffusion policies. *arXiv preprint arXiv:2511.15520*, 2025. 4, 6
- [16] Jianan Li, Xiao Chen, Tao Huang, and Tien-Tsin Wong. Learning to control physically-simulated 3d characters via generating and mimicking 2d motions. *arXiv preprint arXiv:2512.08500*, 2025. 2
- [17] Xinhui Li, Ayush Jain, Zhaojing Yang, Yigit Korkmaz, and Erdem Bryik. When a robot is more capable than a human: Learning from constrained demonstrators. *arXiv preprint arXiv:2510.09096*, 2025. 1, 5
- [18] Yuyang Li, Yinghan Chen, Zihang Zhao, Puhao Li, Tengyu Liu, Siyuan Huang, and Yixin Zhu. Simultaneous tactile-visual perception for learning multimodal robot manipulation. *arXiv preprint arXiv:2512.09851*, 2025. 2, 5
- [19] Ye Li, Jiahe Feng, Yuan Meng, Kangye Ji, Chen Tang, Xinwan Wen, Shutao Xia, Zhi Wang, and Wenwu Zhu. Ts-dp: Reinforcement speculative decoding for temporal adaptive diffusion policy acceleration. *arXiv preprint arXiv:2512.15773*, 2025. 4, 6
- [20] Yifei Li, Haixu Wu, Zeyi Xu, Tuur Stuyck, and Wojciech Matusik. Neural modular physics for elastic simulation. *arXiv preprint arXiv:2512.15083*, 2025. 2, 6
- [21] Zhe Li, Xiang Bai, Jieyu Zhang, Zhuangzhe Wu, Che Xu, Ying Li, Chengkai Hou, and Shanghang Zhang. Urdf-anything: Constructing articulated objects with 3d multimodal language model. *arXiv preprint arXiv:2511.00940*, 2025. 2
- [22] Haotian Liang, Xinyi Chen, Bin Wang, Mingkang Chen, Yitian Liu, Yuhao Zhang, Zanxin Chen, Tianshuo Yang, Yilun Chen, Jiangmiao Pang, Dong Liu, Xiaokang Yang, Yao Mu, Wenqi Shao, and Ping Luo. Mm-act: Learn from multimodal parallel generation to act. *arXiv preprint arXiv:2512.00975*, 2025. 3, 8
- [23] Aileen Liao, Dong-Ki Kim, Max Olan Smith, Ali akbar Agha-mohammadi, and Shayegan Omidshafiei. Delay-aware diffusion policy: Bridging the observation-execution gap in dynamic tasks. *arXiv preprint arXiv:2512.07697*, 2025. 4, 5
- [24] Haowen Liu, Shaoxiong Yao, Haonan Chen, Jiawei Gao, Jiayuan Mao, Jia-Bin Huang, and Yilun Du. Simpackt: Simulation-enabled action planning using vision-language models. *arXiv preprint arXiv:2512.05955*, 2025. 2, 4
- [25] Zhirui Liu, Kaiyang Ji, Ke Yang, Jingyi Yu, Ye Shi, and Jingya Wang. Commanding humanoid by free-form lan-

- 621 guage: A large language action model with unified motion
622 vocabulary. *arXiv preprint arXiv:2511.22963*, 2025. 2
- 623 [26] Patryk Niżeniec and Marcin Iwanowski. Computer vision
624 training dataset generation for robotic environments using
625 gaussian splatting. *arXiv preprint arXiv:2512.13411*, 2025.
626 1, 2, 6
- 627 [27] Jonas Pai, Liam Achenbach, Victoriano Montesinos,
628 Benedek Forrai, Oier Mees, and Elvis Nava. mimic-video:
629 Video-action models for generalizable robot control beyond
630 vlas. *arXiv preprint arXiv:2512.15692*, 2025. 2, 3
- 631 [28] Vineet Pasumarti, Lorenzo Bianchi, and Antonio Loquercio.
632 Agile flight emerges from multi-agent competitive racing.
633 *arXiv preprint arXiv:2512.11781*, 2025. 5
- 634 [29] Markus Peschl, Pietro Mazzaglia, and Daniel Dijkman.
635 From code to action: Hierarchical learning of diffusion-*vlm*
636 policies. *arXiv preprint arXiv:2509.24917*, 2025. 6
- 637 [30] Soujanya Poria, Navonil Majumder, Chia-Yu Hung,
638 Amir Ali Bagherzadeh, Chuan Li, Kenneth Kwok, Zi-
639 wei Wang, Cheston Tan, Jiajun Wu, and David Hsu. 10
640 open challenges steering the future of vision-language-action
641 models. *arXiv preprint arXiv:2511.05936*, 2025. 1, 3, 5
- 642 [31] Massoud Pourmandi. Cooperative perception: A resource-
643 efficient framework for multi-drone 3d scene reconstruc-
644 tion using federated diffusion and nerf. *arXiv preprint*
645 *arXiv:2508.00967*, 2025. 1
- 646 [32] Xiuxiu Qi, Yu Yang, Jiannong Cao, Luyao Bai, Chong-
647 shan Fan, Chengtai Cao, and Hongpeng Wang. Con-
648 tinuous vision-language-action co-learning with semantic-
649 physical alignment for behavioral cloning. *arXiv preprint*
650 *arXiv:2511.14396*, 2025. 5, 8
- 651 [33] Jingjing Qian, Boyao Han, Chen Shi, Lei Xiao, Long Yang,
652 Shaoshuai Shi, and Li Jiang. Geopredict: Leveraging pre-
653 dictive kinematics and 3d gaussian geometry for precise *vla*
654 manipulation. *arXiv preprint arXiv:2512.16811*, 2025. 5
- 655 [34] Yichao Shen, Fangyun Wei, Zhiying Du, Yaobo Liang,
656 Yan Lu, Jiaolong Yang, Nanning Zheng, and Baining Guo.
657 Videovla: Video generators can be generalizable robot ma-
658 nipulators. *The Thirty-ninth Annual Conference on Neural*
659 *Information Processing Systems(NeurIPS2025)*, 2025. 1, 2,
660 6
- 661 [35] Thomas Steinecker, Alexander Bienemann, Denis Trescher,
662 Thorsten Luettel, and Mirko Maehlich. Dynamics-
663 decoupled trajectory alignment for sim-to-real transfer in re-
664 inforcement learning for autonomous driving. *arXiv preprint*
665 *arXiv:2511.07155*, 2025. 1
- 666 [36] Gemini Robotics Team, Coline Devin, Yilun Du, Debidatta
667 Dwivedi, Ruiqi Gao, Abhishek Jindal, Thomas Kipf, Sean
668 Kirmani, Fangchen Liu, Anirudha Majumdar, Andrew Mar-
669 mon, Carolina Parada, Yulia Rubanova, Dhruv Shah, Vikas
670 Sindhvani, Jie Tan, Fei Xia, Ted Xiao, Sherry Yang, Wen-
671 hao Yu, and Allan Zhou. Evaluating gemini robotics policies
672 in a veo world simulator. *arXiv preprint arXiv:2512.10675*,
673 2025. 2, 4
- 674 [37] Nisarg K. Trivedi, Vinayak A. Belludi, Li-Yun Wang, Pardis
675 Taghavi, and Dante Lok. Modest: Multi-optics depth-of-field
676 stereo dataset. *arXiv preprint arXiv:2511.20853*, 2025. 6
- [38] Sümer Tunçay, Alain Andres, and Ignacio Carlucho. Fast
policy learning for 6-dof position control of underwater ve-
hicles. *arXiv preprint arXiv:2512.13359*, 2025. 5
- [39] Federico Vasile, Ri-Zhao Qiu, Lorenzo Natale, and Xiao-
long Wang. Gaussian-augmented physics simulation and
system identification with complex colliders. *arXiv preprint*
arXiv:2511.06846, 2025. 2
- [40] Hanshi Wang, Yuhao Xu, Zekun Xu, Jin Gao, Yufan Liu,
Weiming Hu, Ke Wang, and Zhipeng Zhang. Autoprune:
Each complexity deserves a pruning policy. *NeurIPS 2025*,
2025. 4, 6
- [41] Maggie Wang, Stephen Tian, Aiden Swann, Ola Shorinwa,
Jiajun Wu, and Mac Schwager. Phys2real: Fusing *vlm* pri-
ors with interactive online adaptation for uncertainty-aware
sim-to-real manipulation. *arXiv preprint arXiv:2510.11689*,
2025. 1, 4
- [42] Sen Wang, Jingyi Tian, Le Wang, Zhimin Liao, Jiayi Li,
Huaiyi Dong, Kun Xia, Sanping Zhou, Wei Tang, and
Hua Gang. Sampo:scale-wise autoregression with mo-
tion prompt for generative world models. *arXiv preprint*
arXiv:2509.15536, 2025. 1, 3
- [43] Weitian Wang, Lukas Meiner, Rai Shubham, Cecilia De La
Parra, and Akash Kumar. Htm: Head-wise temporal token
merging for faster *vggt*. *arXiv preprint arXiv:2511.21317*,
2025. 2
- [44] Hao Wu, Yuan Gao, Xingjian Shi, Shuaipeng Li, Fan Xu,
Fan Zhang, Zhihong Zhu, Weiyang Wang, Xiao Luo, Kun
Wang, Xian Wu, and Xiaomeng Huang. Spatiotemporal
forecasting as planning: A model-based reinforcement learn-
ing approach with generative world models. *arXiv preprint*
arXiv:2510.04020, 2025. 1, 4, 6
- [45] Xiaoshan Wu, Yifei Yu, Xiaoyang Lyu, Yihua Huang,
Bo Wang, Baoheng Zhang, Zhongrui Wang, and Xiao-
juan Qi. Eag3r: Event-augmented 3d geometry estimation
for dynamic and extreme-lighting scenes. *arXiv preprint*
arXiv:2512.00771, 2025. 1, 4
- [46] Zixuan Wu, Hengyuan Zhang, Ting-Hsuan Chen, Yuliang
Guo, David Paz, Xinyu Huang, and Liu Ren. Dino-diffusion
modular designs bridge the cross-domain gap in autonomous
parking. *arXiv preprint arXiv:2510.20335*, 2025. 3
- [47] Tianyu Xu, Jiawei Chen, Jiazhao Zhang, Wenyao Zhang,
Zekun Qi, Minghan Li, Zhizheng Zhang, and He Wang.
Mm-nav: Multi-view *vla* model for robust visual navigation
via multi-expert learning. *arXiv preprint arXiv:2510.03142*,
2025. 5
- [48] Wenjiang Xu, Cindy Wang, Rui Fang, Mingkan Zhang,
Lusong Li, Jing Xu, Jiayuan Gu, Zecui Zeng, and Rui
Chen. Embodied tree of thoughts: Deliberate manipula-
tion planning with embodied world model. *arXiv preprint*
arXiv:2512.08188, 2025. 3, 4
- [49] Xinyue Xu, Jieqiang Sun, Jing, Dai, Siyuan Chen, Lanjie
Ma, Ke Sun, Bin Zhao, Jianbo Yuan, Sheng Yi, Haohua Zhu,
and Yiwen Lu. Dexcanvas: Bridging human demonstra-
tions and robot learning for dexterous manipulation. *arXiv*
preprint arXiv:2510.15786, 2025. 2, 5
- [50] Mingxuan Yan, Yuping Wang, Zechun Liu, and Jiachen
Li. Rdd: Retrieval-based demonstration decomposer for

- 734 planner alignment in long-horizon tasks. *arXiv preprint*
735 *arXiv:2510.14968*, 2025. 4, 5
- 736 [51] Qianfeng Yang, Xiang Chen, Pengpeng Li, Qiyuan Guan,
737 Guiyue Jin, and Jiyu Jin. Rethinking rainy 3d scene recon-
738 struction via perspective transforming and brightness tuning.
739 *arXiv preprint arXiv:2511.06734*, 2025. 1, 6
- 740 [52] Jianglong Ye, Lai Wei, Guangqi Jiang, Changwei Jing,
741 Xueyan Zou, and Xiaolong Wang. From power to precision:
742 Learning fine-grained dexterity for multi-fingered robotic
743 hands. *arXiv preprint arXiv:2511.13710*, 2025. 5
- 744 [53] Maryam Yousefi and Soodeh Bakhshandeh. Curvature-
745 regularized variational autoencoder for 3d scene reconstruc-
746 tion from sparse depth. *arXiv preprint arXiv:2512.05783*,
747 2025. 1, 4
- 748 [54] Ziyao Zeng, Jingcheng Ni, Ruyi Liu, and Alex Wong.
749 Coffee: Controllable diffusion fine-tuning. *arXiv preprint*
750 *arXiv:2511.14113*, 2025. 2, 3, 6
- 751 [55] Kevin Zhang, Kuangzhi Ge, Xiaowei Chi, Renrui Zhang,
752 Shaojun Shi, Zhen Dong, Sirui Han, and Shanghang Zhang.
753 Can world models benefit vlms for world dynamics? *arXiv*
754 *preprint arXiv:2510.00855*, 2025. 2, 3
- 755 [56] Junwei Zhou and Yu-Wing Tai. Amodalgen3d: Generative
756 amodal 3d object reconstruction from sparse unposed views.
757 *arXiv preprint arXiv:2511.21945*, 2025. 1, 6
- 758 [57] Jian Zhou, Sihao Lin, Shuai Fu, and Qi WU. Decoupled ac-
759 tion head: Confining task knowledge to conditioning layers.
760 *arXiv preprint arXiv:2511.12101*, 2025. 4, 6
- 761 [58] Sizhuo Zhou, Xiaosong Jia, Fanrui Zhang, Junjie Li, Juyong
762 Zhang, Yukang Feng, Jianwen Sun, Songbur Wong, Junqi
763 You, and Junchi Yan. Lagen: Towards autoregressive lidar
764 scene generation. *arXiv preprint arXiv:2511.21256*, 2025. 1,
765 3
- 766 [59] Minghan Zhu, Zhiyi Wang, Qihang Sun, Maani Ghaffari,
767 and Michael Posa. Object reconstruction under occlusion
768 with generative priors and contact-induced constraints. *arXiv*
769 *preprint arXiv:2512.05079*, 2025. 1, 6
- 770 [60] Ruijie Zhu, Mulin Yu, Linning Xu, Lihan Jiang, Yixuan Li,
771 Tianzhu Zhang, Jiangmiao Pang, and Bo Dai. Objectgs:
772 Object-aware scene reconstruction and scene understanding
773 via gaussian splatting. *arXiv preprint arXiv:2507.15454*,
774 2025. 1, 2