

InterPrior: Scaling Generative Control for Physics-Based Human-Object Interactions

Anonymous CVPR submission

Paper ID ****

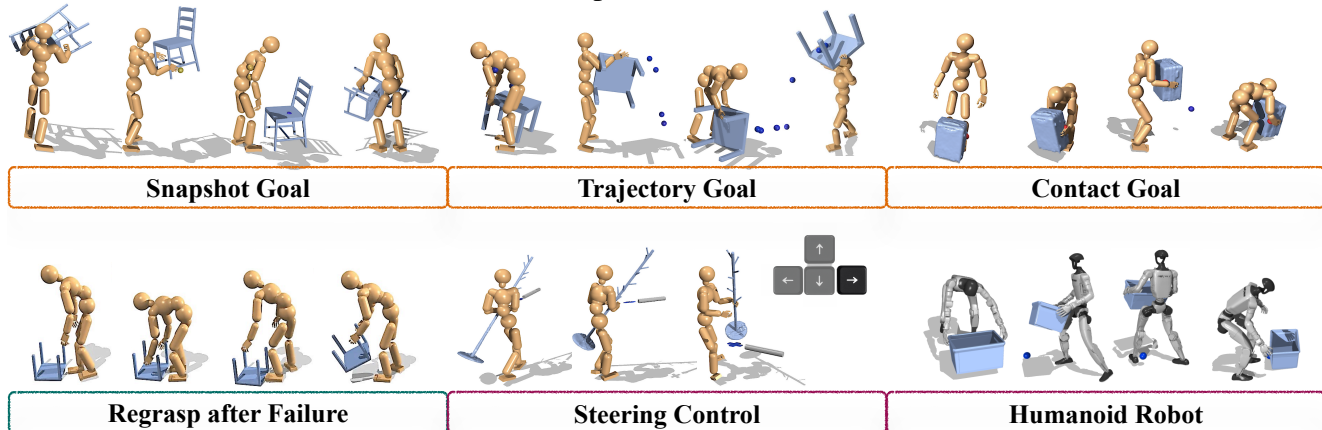


Figure 1. InterPrior is a versatile generative controller that drives a simulated humanoid to follow sparse goals and interact with objects in physics. It supports (I) long-horizon snapshot goals, (II) trajectory goals, and (III) contact goals (Top). Yellow, blue, and red dots respectively denote human, object, and contact goals. It demonstrates failure recovery (Bottom Left) and steering control on humanoid robot embodiments (Bottom Right).

Abstract

001 Humans rarely plan whole-body interactions with objects
 002 at the level of explicit whole-body movements. High-level
 003 intentions, such as affordance, define the goal, while co-
 004 ordinated balance, contact, and manipulation can emerge
 005 naturally from underlying physical and motor priors. Scal-
 006 ing such priors is key to enabling humanoids to compose
 007 and generalize loco-manipulation skills across diverse con-
 008 texts while maintaining physically coherent whole-body co-
 009 ordination. To this end, we introduce InterPrior, a scal-
 010 able framework that learns a unified generative controller
 011 through large-scale imitation pretraining and post-training
 012 by reinforcement learning. InterPrior first distills a full-
 013 reference imitation expert into a versatile, goal-conditioned
 014 variational policy that reconstructs motion from multimodal
 015 observations and high-level intent. While the distilled pol-
 016 icy reconstructs training behaviors, it does not generalize
 017 reliably due to the vast configuration space of large-scale
 018 human-object interactions. To address this, we apply data
 019 augmentation with physical perturbations, and then per-
 020 form reinforcement learning finetuning to improve compe-
 021 tence on unseen goals and initializations. Together, these
 022 steps consolidate the reconstructed latent skills into a valid

manifold, yielding a motion prior that generalizes beyond 023
 the training data, e.g., it can incorporate new behaviors 024
 such as interactions with unseen objects. We further demon- 025
 strate its effectiveness for user-interactive control and its 026
 potential for real robot deployment. 027

1. Introduction 028

Human-object interaction (HOI) is hierarchical: humans 029
 plan with sparse intent while limb coordination, balance, 030
 and contact emerge through fast motor responses [16]. 031
 Imitation policies [21] scale to large HOI skill repertoires but 032
 rely on explicit planners providing dense full-body and ob- 033
 ject references. A more useful interaction motor prior 034
 should sample feasible loco-manipulation behaviors from 035
 a distribution conditioned on sparse goals, rather than mim- 036
 icking deterministic, fully-specified trajectories. 037

Existing approaches struggle to scale. Generative con- 038
 trollers trained with adversarial matching [3, 10] suffer 039
 from unstable optimization and handcrafted task rewards. 040
 Reference-imitation distillation [8, 14] absorbs large-scale 041
 data but is brittle when reference coverage lags the config- 042
 uration space, as in loco-manipulation, where a few object 043
 DoFs combinatorially explode contact modes. 044

We introduce InterPrior, a generative HOI controller 045

(Fig. 1) scalable along four axes: *task coverage* (one policy supports sparse targets and their compositions), *skill coverage* (a single recipe scales to large HOI data), *motion coverage* (it generates expressive trajectories beyond demonstrations), and *dynamics coverage* (it sustains success under varied physics). Our key insight is that *RL finetuning is essential* for turning distillation from data reconstruction into a robust, generalizable policy. Distillation alone cannot cover the full HOI configuration space, while RL alone drifts toward unnatural reward-hacking. We use distillation as a strong, natural initialization, and RL as a local optimizer that improves robustness while remaining anchored to the pretrained model.

2. Related Work

Physics-based HOI control. Reference imitation policies [19, 21] produce high-fidelity HOI motion when dense references are available, but rely on explicit kinematic planners that scale poorly with object DoFs. Adversarial generative controllers [3, 10] broaden distributions but suffer from unstable optimization and discriminator mode collapse. Goal-conditioned distillation [15] absorbs large-scale data without task-specific design, yet remains brittle when reference coverage lags the contact-rich configuration space. InterPrior combines imitation distillation with RL post-training to consolidate sparse-goal coverage and recovery in a single policy.

Humanoid loco-manipulation. Whole-body humanoid policies have shown locomotion and pick-and-place [2, 5], often via teleoperation or staged imitation. Recent works combine retargeting with universal trackers for general loco-manipulation [4]; InterPrior acts as the underlying interaction prior such trackers can build on.

3. Method

Task formulation. We learn a policy π that operates in a physics simulator and produces HOI motion from *high-level goals* rather than full reference. Goals can be supplied by a human user, a kinematic motion generator, or sparse keypoints from MoCap. The policy conditions on the current human-object state and recent history together with these goals, and samples control signals from its learned distribution. The output is a physically simulated rollout that follows the goals where specified and remains diverse elsewhere.

Goal conditioning. At each timestep the policy receives short-horizon previews and a long-horizon snapshot, each represented as a *masked residual encoding*: a binary mask selects which goal components (joints, object pose, contacts) are active, and active entries are encoded as residuals to the current state. This unifies snapshot, trajectory, contact, and composed goals under a single masked formu-

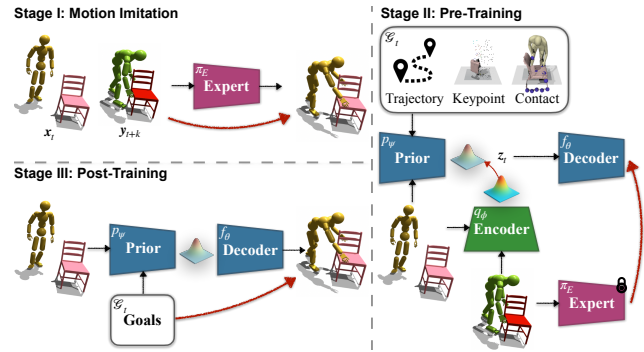


Figure 2. Overview of InterPrior. (I) Train an imitation expert on large-scale HOI data with augmentation, perturbations, and a contact-aware hand reward; (II) distill the expert into a variational policy with structured latent skills via online DAGger and ELBO loss; (III) RL post-training on randomized single-frame in-betweening, anchored by concurrent distillation. Blue: inference-time modules. Green/red: training-only.

lation [14]. At inference, user-specified or model-generated sparse targets are supplied by filling only the informed components and zeroing the rest.

Stage I: imitation expert. We train an expert π_E on large-scale HOI data with PPO [13]. Beyond reference tracking, we expand reference scope via random initialization, sparse impulse perturbations, and dynamics randomization, and we introduce a *reference-free hand reward* that encourages the hand to target and wrap the actual simulated object rather than rigidly tracking the reference. This corrects failures on thin geometries and under perturbation, where strict reference tracking becomes unreliable.

Stage II: variational distillation. We distill π_E into a masked conditional variational policy π with a Transformer prior, an MLP encoder (training-only), and an MLP decoder. Latents are projected onto a hypersphere ($z_t \leftarrow z_t / \|z_t\|$) following [11] to bound out-of-distribution draws while preserving directional diversity. Online DAGger [12] minimizes a weighted ELBO [6] (imitation, masked-goal reconstruction, KL), with two auxiliary losses constraining the prior mean to unit magnitude and encouraging temporal consistency across consecutive priors. The decoder also reconstructs masked goal entries, yielding a latent that learns to *complete* intent.

Stage III: RL post-training. The distilled policy follows goals but is brittle when state or goals drift off the demonstration manifold. We RL-finetune on *single-frame in-betweening*: from randomized initial states, the policy must reach a single sparse goal sampled from the dataset, with a sparse success reward $r_{\text{goal}} = r_{\text{succ}} \cdot \mathbb{1}[\|\mathbf{m}_{t+L} \odot \Delta(\hat{\mathbf{y}}_{t+L}, \mathbf{x}_t)\|_1 < \tau]$. Composed and randomized goals systematically broaden the encountered state distribution, training the policy to recover from near-failure states (e.g., regrasping after a slipped grasp) without explicit supervision. To preserve the pretrained prior, a subset of environ-

Table 1. **Goal-conditioned tasks** (% success). Snapshot/Trajectory/Contact in-distribution tasks plus Chain and Random-Init stress tests; each row adds one component.

Variant	Snap	Traj	Cont	Chain	RInit	Fail ↓
MaskedMimic [14]	64.2	88.0	52.2	29.1	31.7	12.6
+ InterMimic+ Expert	71.4	92.7	69.3	33.9	30.1	11.0
+ Latent Shaping	74.9	92.4	71.9	40.0	30.9	10.6
+ Bounded Latent & Obs.	89.1	93.6	88.5	45.1	41.1	6.0
+ RL Finetuning (full)	90.0	94.6	90.7	68.8	88.6	3.7

Table 2. **Full-reference tracking and adaptation.** OMOMO uses thin objects with initialization perturbations; BEHAVE/HODome report zero-shot transfer to unseen objects.

Method	OMOMO SR ↑	BEHAVE SR ↑	HODome SR ↑
InterMimic [21]	63.9	10.7	27.8
InterMimic + finetuning	–	38.9	55.5
InterPrior (Ours)	83.2	27.4	40.1
InterPrior + finetuning	–	52.0	72.4

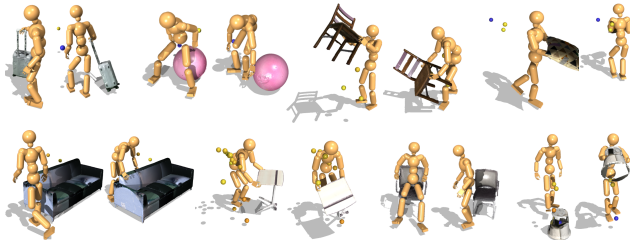


Figure 3. **Zero-shot generalization.** A single InterPrior model trained on OMOMO [7] adapts to unseen objects and interactions from BEHAVE [1] and HODome [22].

131 ments continues optimizing the distillation objective in parallel, anchoring the policy to its natural skill manifold without freezing weights. Skills outside demonstrations (e.g., 132 *getting up*) are introduced via a learnable subtask token and 133 a shaped upright-posture reward. 134 135

136 4. Experiments

137 **Setup.** We train on the InterAct [20] preprocessing of 138 OMOMO [7] and evaluate generalization on selected sequences 139 from BEHAVE [1] and HODome [22]. Policies run at 30 Hz in IsaacGym [9]. We compare InterPrior 140 against the original InterMimic [21] (full-reference imitator) 141 and an adapted MaskedMimic [14, 15] (sparse-goal distillation). 142 Tasks include (i) full-reference tracking on thin objects with 143 initialization noise; (ii) sparse goal following under snapshot, 144 trajectory, contact, and composed goals; (iii) stress tests with 145 long-horizon multi-goal chains and random-init lifting; and 146 (iv) zero-shot adaptation to unseen objects/interactions. 147 148

149 **Goal-conditioned results.** Table 1 shows cumulative gains 150 as we add each component. The largest jumps appear 151 on stress tests: bounded latent and observation spaces 152 roughly double random-init success, and RL finetuning further 153 raises it from 41.1% to 88.6% while halving failure rate

154 to 3.7%. Distillation-only policies fit the demonstration- 155 induced state distribution; long rollouts and goal switching 156 enter under-covered states and drift, while RL on single- 157 frame in-betweening directly trains for sparse-target reach- 158 ing from diverse initializations, improving recovery from 159 off-distribution states. Trajectory-following is preserved 160 because finetuning operates on snapshot goals while concurrent 161 distillation protects trajectory-conditioned modes. 162

163 **Tracking and adaptation.** Table 2 shows that InterPrior 164 achieves **83.2%** success on OMOMO with thin objects and 165 initialization noise, versus 63.9% for InterMimic, while remaining 166 competitive on per-joint position errors. The gain stems from 167 the contact-aware hand reward and bounded latent: the policy 168 makes small targeted deviations to realign contact instead of 169 rigidly tracking the reference. On BEHAVE and HODome (held- 170 out objects, different human shapes), InterPrior and its finetuned 171 variant outperform InterMimic by large margins, showing the 172 distilled prior absorbs additional interaction data more effectively 173 even when that data is imperfect. 174

175 **Qualitative behavior.** Fig. 1 shows minute-long whole-body 176 interaction with smooth transitions across approach, grasp, lift, 177 and reposition; when contact or balance drift, InterPrior self- 178 corrects rather than compounding errors. Fig. 3 shows zero-shot 179 generalization: sparse snapshot goals are sufficient for InterPrior 180 to complete unspecified DoFs and converge to feasible contacts 181 on unseen objects. 182

183 **Embodiment and sim-to-real.** We retrain InterPrior on 184 the Unitree G1 [18] with G1-specific stabilization rewards and 185 dynamics randomization, and observe sustained loco-manipulation 186 under sim-to-sim transfer to MuJoCo [17]. Building on this prior, 187 a downstream controller [4] deploys autonomous loco-manipulation 188 on a real G1 from egocentric depth and sparse task goals, indicating 189 that scaling interaction priors in simulation is a viable path to 190 versatile real-world humanoid behavior. 191

192 **Operator-controlled steering.** Beyond dataset-derived 193 goals, InterPrior accepts on-the-fly sparse targets from a human 194 operator via a keyboard interface that streams a small number 195 of object-pose or end-effector commands per second. Because the 196 policy is conditioned on the same masked goal interface used at 197 training, no additional finetuning is needed to support this mode. 198 In our deployments, an operator chains push, lift, and reposition 199 commands across multiple objects, and the policy completes the 200 unspecified DoFs while maintaining whole-body balance. Mid-trajectory 201 command switches are absorbed gracefully, and brief contact loss 202 triggers re-approach behaviors learned through randomized in- 203 betweening rather than scripted primitives. 204

205 **Implementation details.** The imitation expert, encoder, and 206 decoder are MLPs with hidden sizes (1024, 1024, 512); the prior 207 is a 4-layer Transformer encoder; critics share the expert MLP 208 architecture. Latent dimension is $d_z = 64$. 209

207 All policies are trained with PPO [13]; distillation uses on-
208 line DAgger. For the G1 humanoid retraining, we addition-
209 ally randomize friction, mass density, center-of-mass off-
210 sets, and inertia, and apply per-joint impulse perturbations
211 during rollout. Episode-fixed sampling noise $\epsilon \sim \mathcal{N}(0, I)$
212 promotes temporal consistency across a rollout.

213 **Failure modes.** Remaining failures involve extremely thin
214 or elongated unseen geometries, and partial completion in
215 multi-goal chaining when canonicalization across subgoals
216 introduces large alignment discrepancies, where the pol-
217 icy favors balance over precise goal matching. Integrat-
218 ing tactile sensing or depth-conditioned goal completion is
219 a promising direction.

220 5. Conclusion

221 InterPrior is a generative HOI controller that combines
222 large-scale imitation distillation with RL post-training un-
223 der perturbations and randomized goals. Distillation pro-
224 vides a strong, natural initialization; RL anchors the policy
225 to its pretrained skill manifold while expanding state cov-
226 erage and recovery behavior. The result is a reusable inter-
227 action prior that handles snapshot, trajectory, contact, and
228 composed goals; recovers from failures; generalizes zero-
229 shot to unseen objects; and transfers across embodiments
230 toward real-humanoid deployment.

231 References

232 [1] Bharat Lal Bhatnagar, Xianghui Xie, Ilya Petrov, Cristian
233 Sminchisescu, Christian Theobalt, and Gerard Pons-Moll.
234 BEHAVE: Dataset and method for tracking human object in-
235 teractions. In *CVPR*, 2022. 3

236 [2] Yuhui Fu, Feiyang Xie, Chaoyi Xu, Jing Xiong, Haoqi Yuan,
237 and Zongqing Lu. DemoHLM: From one demonstration to
238 generalizable humanoid loco-manipulation. *arXiv preprint*
239 *arXiv:2510.11258*, 2025. 2

240 [3] Mohamed Hassan, Yunrong Guo, Tingwu Wang, Michael
241 Black, Sanja Fidler, and Xue Bin Peng. Synthesizing phys-
242 ical character-scene interactions. In *SIGGRAPH*, 2023. 1,
243 2

244 [4] Xialin He, Sirui Xu, Xinyao Li, Runpei Dong, Liuyu Bian,
245 Yu-Xiong Wang, and Liang-Yan Gui. ULTRA: Unified mul-
246 timodal control for autonomous humanoid whole-body loco-
247 manipulation. *arXiv preprint arXiv:2603.03279*, 2026. 2,
248 3

249 [5] Dvij Kalaria, Sudarshan S Harithas, Pushkal Katara,
250 Sangkyung Kwak, Sarthak Bhagat, Shankar Sastry, Srinath
251 Sridhar, Sai Vemprala, Ashish Kapoor, and Jonathan Chung-
252 Kuan Huang. DreamControl: Human-inspired whole-body
253 humanoid control for scene interaction via guided diffusion.
254 *arXiv preprint arXiv:2509.14353*, 2025. 2

255 [6] Diederik P Kingma and Max Welling. Auto-encoding varia-
256 tional bayes. *arXiv preprint arXiv:1312.6114*, 2013. 2

257 [7] Jiaman Li, Jiajun Wu, and C Karen Liu. Object motion
258 guided human motion synthesis. *ACM Transactions on*
259 *Graphics (TOG)*, 42(6):1–11, 2023. 3

[8] Zhengyi Luo, Jinkun Cao, Sammy Christen, Alexander Win-
260 kler, Kris Kitani, and Weipeng Xu. Grasping diverse objects
261 with simulated humanoids. In *NeurIPS*, 2024. 1
262

[9] Viktor Makoviychuk, Lukasz Wawrzyniak, Yunrong Guo,
263 Michelle Lu, Kier Storey, Miles Macklin, David Hoeller,
264 Nikita Rudin, Arthur Allshire, Ankur Handa, et al. Isaac
265 gym: High performance gpu-based physics simulation for
266 robot learning. In *NeurIPS*, 2021. 3
267

[10] Liang Pan, Zeshi Yang, Zhiyang Dou, Wenjia Wang, Buzhen
268 Huang, Bo Dai, Taku Komura, and Jingbo Wang. Token-
269 HSI: Unified synthesis of physical human-scene interac-
270 tions through task tokenization. In *CVPR*, 2025. 1, 2
271

[11] Xue Bin Peng, Yunrong Guo, Lina Halper, Sergey Levine,
272 and Sanja Fidler. Ase: Large-scale reusable adversarial
273 skill embeddings for physically simulated characters. *ACM*
274 *Transactions On Graphics (TOG)*, 41(4):1–17, 2022. 2
275

[12] Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A re-
276 duction of imitation learning and structured prediction to no-
277 regret online learning. In *Proceedings of the fourteenth inter-
278 national conference on artificial intelligence and statistics*,
279 pages 627–635. JMLR Workshop and Conference Proceed-
280 ings, 2011. 2
281

[13] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Rad-
282 ford, and Oleg Klimov. Proximal policy optimization algo-
283 rithms. *arXiv preprint arXiv:1707.06347*, 2017. 2, 4
284

[14] Chen Tessler, Yunrong Guo, Ofir Nabati, Gal Chechik, and
285 Xue Bin Peng. Maskedmimic: Unified physics-based char-
286 acter control through masked motion inpainting. *ACM Trans-
287 actions on Graphics (TOG)*, 43(6):1–21, 2024. 1, 2, 3
288

[15] Chen Tessler, Yifeng Jiang, Erwin Coumans, Zhengyi Luo,
289 Gal Chechik, and Xue Bin Peng. MaskedManipulator:
290 Versatile whole-body control for loco-manipulation. *arXiv*
291 *preprint arXiv:2505.19086*, 2025. 2, 3
292

[16] Emanuel Todorov and Michael I Jordan. Optimal feedback
293 control as a theory of motor coordination. *Nature neuro-
294 science*, 5(11):1226–1235, 2002. 1
295

[17] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A
296 physics engine for model-based control. In *IROS*, 2012. 3
297

[18] Unitree. Unitree g1 humanoid agent ai avatar. [https://
298 www.unitree.com/g1/](https://www.unitree.com/g1/). 3
299

[19] Yinhuai Wang, Jing Lin, Ailing Zeng, Zhengyi Luo, Jian
300 Zhang, and Lei Zhang. PhysHOI: Physics-based imita-
301 tion of dynamic human-object interaction. *arXiv preprint*
302 *arXiv:2312.04393*, 2023. 2
303

[20] Sirui Xu, Dongting Li, Yucheng Zhang, Xiyan Xu, Qi Long,
304 Ziyin Wang, Yunzhi Lu, Shuchang Dong, Hezi Jiang, Ak-
305 shat Gupta, Yu-Xiong Wang, and Liang-Yan Gui. InterAct:
306 Advancing large-scale versatile 3d human-object interaction
307 generation. In *CVPR*, 2025. 3
308

[21] Sirui Xu, Hung Yu Ling, Yu-Xiong Wang, and Liang-Yan
309 Gui. InterMimic: Towards universal whole-body control for
310 physics-based human-object interactions. In *CVPR*, 2025. 1,
311 2, 3
312

[22] Juze Zhang, Haimin Luo, Hongdi Yang, Xinru Xu, Qianyang
313 Wu, Ye Shi, Jingyi Yu, Lan Xu, and Jingya Wang. Neural-
314 Dome: A neural modeling pipeline on multi-view human-
315 object interactions. In *CVPR*, 2023. 3
316