# Human2LocoMan: Learning Versatile Quadrupedal Manipulation with Human Pretraining

Yaru Niu<sup>1\*</sup>, Yunzhe Zhang<sup>1\*</sup>, Mingyang Yu<sup>1</sup>, Changyi Lin<sup>1</sup>, Chenhao Li<sup>1</sup>, Yikai Wang<sup>1</sup>, Yuxiang Yang<sup>2</sup>, Wenhao Yu<sup>2</sup>, Tingnan Zhang<sup>2</sup>, Zhenzhen Li<sup>3</sup>, Jonathan Francis<sup>1,3</sup>, Bingqing Chen<sup>3</sup>, Jie Tan<sup>2</sup>, and Ding Zhao<sup>1</sup> <sup>1</sup>Carnegie Mellon University <sup>2</sup>Google DeepMind <sup>3</sup>Bosch Center for Artificial Intelligence \*Equal contributions



Fig. 1: Human2LocoMan provides a unified framework for collecting human demonstrations and teleoperated robot wholebody motions, enabling flexible and scalable data collection. Human data is used for cross-embodiment model pretraining, while robot data is leveraged for policy finetuning. Human2LocoMan achieves positive transfer from human to quadrupedal embodiments, facilitating versatile quadrupedal manipulation.

Abstract-Ouadrupedal robots have demonstrated impressive locomotion capabilities in complex environments, but equipping them with autonomous versatile manipulation skills in a scalable way remains a significant challenge. In this work, we introduce a system that integrates data collection and imitation learning from both humans and LocoMan, a quadrupedal robot with multiple manipulation modes. Specifically, we introduce a teleoperation and data collection pipeline, supported by dedicated hardware, which unifies and modularizes the observation and action spaces of the human and the robot. To effectively leverage the collected data, we propose an efficient learning architecture that supports co-training and pretraining with multimodal data across different embodiments. Additionally, we construct the first manipulation dataset for the LocoMan robot, covering various household tasks in both unimanual and bimanual modes, supplemented by a corresponding human dataset. Experimental results demonstrate that our data collection and training framework significantly improves the efficiency and effectiveness of imitation learning, enabling more versatile quadrupedal manipulation capabilities. Our hardware, data, and code will be open-sourced at: https://human2bots.github.io.

# I. INTRODUCTION

While quadrupedal robots have demonstrated impressive locomotion capabilities in complex environments [1]–[7],

and recent advances have extended their abilities to manipulation tasks [8]-[14], enabling autonomous and versatile quadrupedal manipulation at scale remains a major challenge. In this work, we take inspiration from the open-source LocoMan platform [14], a quadrupedal robot equipped with two leg-mounted loco-manipulators, which offers a versatile foundation for learning manipulation skills across multiple operating modes. Imitation learning has long been a fundamental approach for teaching robots complex skills through demonstrations [15], with the acquisition of high-quality data being critical for achieving efficient and effective learning. Prior works have explored various strategies for collecting in-domain robot data, primarily focusing on robot arms [16]-[19], humanoid robots [20]-[22], and quadrupeds equipped with top-mounted arms [10], [11], [23]. However, collecting egocentric manipulation data on a quadrupedal platform like LocoMan remains underexplored. To scale up data collection for imitation learning, recent works propose leveraging simulation data [24]–[26] or human data [17], [27]–[31]. Human data, in particular, have been used to provide highlevel task guidance [17], [28], improve visual encoders [29], simulate in-domain robot data [27], [30], or serve as additional training data by treating humans as an alternative



Fig. 2: Human2LocoMan framework. (a) The data collection system leverages an XR headset to collect egocentric human data and teleoperated robot data. Human and robot data are mapped to a unified coordinate frame. (b) The dataset consists of aligned vision, proprioception, and actions from the human and the robot. (c) During training, the network is first pretrained on easy-to-collect human data, and then finetuned on a small amount of robot data. (d) We evaluate the autonomous Human2LocoMan policies on six household tasks in unimanual and bimanual modes.

embodiment with similar kinematic structures [31]. However, transferring skills from humans to quadrupedal robots remains challenging due to the substantial embodiment gap, which complicates both data collection and policy transfer. To address these challenges, we propose Human2LocoMan, a unified framework that bridges the human-to-quadruped gap. Human2LocoMan introduces a novel teleoperation and data collection system that aligns human and robot data, coupled with a modular transformer-based architecture for robust cross-embodiment learning. Together, these components enable scalable learning of versatile manipulation skills on quadrupedal robots.

Specifically, to enable scalable data collection, our system leverages an extended reality (XR) headset to capture human motions while streaming a first-person or first-robot (during teleoperation) view to the operator. For human data collection, the operator simply wears the XR headset and performs tasks naturally. During teleoperation, we align the human and quadruped into a unified coordinate frame to bridge the embodiment gap. In addition to mapping human hand motions to the robot's grippers, we map human head motions to the robot's torso, expanding the robot's workspace and enhancing active sensing capabilities. Target poses are then passed to a whole-body controller to generate coordinated robot motions.

In contrast to works that use egocentric human data to pretrain vision encoders [29] or learn high-level intent [17], we treat the human as another embodiment and use human data for cross-embodiment learning. Despite mapping human and robot data to a unified frame, there exist obvious gaps ranging from differences in dynamics to extra wrist cameras on the robot. Thus, we design a modular transformer architecture, *Modularized Cross-embodiment Transformer* (MXT), which shares the transformer trunk, but has embodiment-specific tokenizers / detokenizers. To enable positive transfer, the MXT policy is first pretrained on human data and subsequently finetuned with a small amount of robot data. We evaluate our approach on six household tasks, across both unimanual and bimanual manipulation modes. Our results demonstrate strong task performance by MXT compared to competitive baselines, effective positive transfer from human demonstrations to robot policies, and increased robustness to both in-distribution (ID) and out-of-distribution (OOD) scenarios.

# II. METHODOLOGY

In this section, we present the design and implementation of our system, Human2LocoMan, which integrates teleoperation, data collection, and a Transformer-based architecture for cross-embodied learning.

## A. Human2LocoMan System Overview

We utilize the Apple Vision Pro headset and the Open-Television system [21] to capture human motions and stream first-person or first-robot video to the human operator. A lightweight stereo camera with a 120-degree horizontal field of view is mounted on both the VR headset and the LocoMan robot to provide egocentric views, while additional cameras, such as RGB wrist cameras, can be optionally attached to the robot. Through the Human2LocoMan teleoperation system (Section II-B), the human operator can control the LocoMan robot to perform versatile manipulation tasks in both unimanual and bimanual modes. In the unimanual mode, we also map human head motions to the robot's torso movements to expand the teleoperation workspace and enhance active sensing. The Human2LocoMan system enables the collection of both human and robot data, transforming them into a shared space. Masks are applied to distinguish across different embodiments and manipulation modes. The collected human data are used to pretrain an action model called the Modularized Cross-embodiment Transformer (MXT). The in-domain robotic data collected via teleoperation are used to finetune the pretrained model to learn a manipulation policy that predicts the 6D poses of LocoMan's end effectors and torso, as well as gripper actions.



Fig. 3: **Modularized Cross-embodiment Transformer (MXT) architecture.** The inputs are organized as a list of modalities and encoded each by a separate tokenizer into a fixed number of tokens. The transformer trunk handles decision making by consuming the concatenated encoded tokens and producing a fixed number of raw output tokens. Each of the detokenizers at the end decodes a fixed subset of the output tokens into a modality of the final actions.

#### B. Human2LocoMan Teleoperation and Data Collection

A unified frame for both human and LocoMan. To map human motions to LocoMan's various operation modes via VR-based teleoperation—and to enhance the transferability of motion data across different embodiments—we establish a unified reference frame,  $\mathcal{F}_u$ , to align motions across embodiments. As shown in Figure 2 (a), this unified frame is attached to the rigid body where the main camera is mounted. At the embodiment's reset pose, the x-axis points forward, aligned with the workspace and parallel to the ground; the y-axis points leftward; and the z-axis points upward, perpendicular to the ground.

Motion mapping. We map the human wrist motions to LocoMan's end-effector motions, map the human head motions to LocoMan's torso motions, and hand poses to LocoMan's gripper actions. The 6D poses of the human hand, head, and wrist poses in SE(3) in the VR-defined world frame are streamed from the VR set to the Human2LocoMan teleoperation server. Please refer to IV-A for more details.

**Whole-body controller.** The robot target pose at time t,  $p_t^t$ , is calculated from the teleoperation server, and sent to the whole-body controller of the LocoMan robot, which is adapted from the one introduced in [14]. Please refer to Appendix Section IV-B for more details.

**Data Collection.** The details of Human2LocoMan data collection can be found in Appendix Section IV-C. We ensure that human and robot data are unified in both format and spatial interpretation, and can be used to train our proposed Modularized Cross-Embodiment Transformer introduced in Section II-C.

# C. Modularized Cross-embodiment Transformer

To train a policy on LocoMan that benefits from heterogeneous human data, we opt for task-space control in this work, where the actions predicted by the policy are represented as key pose parameters of the physical embodiment, such as the end effector 6D pose and the body 6D pose. Given our unified multi-embodiment data collection pipeline, we aim to train a cross-embodiment policy where the overall structure and the majority of parameters are transferrable. To this end, we propose a modularized design called *Modularized Cross*- *embodiment Transformer* (MXT). MXT consists mainly of three groups of modules: tokenizers, transformer trunk, and detokenizers. The tokenizers act as encoders and map embodiment-specific observations to tokens in the latent space, and the detokenizers translate the output tokens from the trunk to actions in the action space of each embodiment. The tokenizers and detokenizers are specific to one embodiment and are reinitialized for each new embodiment, while the trunk is shared across all embodiments and reused for transferring the policy among embodiments. Figure 3 illustrates the architecture of our network. The design details and training paradigm of MXT are elaborated in Appendix Section IV-D and IV-E, respectively.

## III. EXPERIMENTS

In this section, we aim to answer the following research questions: (1) Does the Human2LocoMan system enable versatile quadrupedal manipulation capabilities? (2) How does MXT compare to state-of-the-art imitation learning architectures? (3) How does human data collected by Human2LocoMan contribute to imitation learning performance? (4) Do the design choices in MXT facilitate positive transfer from Human to LocoMan?

#### A. Experimental Setup

We evaluate MXT on six diverse household manipulation tasks—unimanual and bimanual toy collection, unimanual and bimanual shoe rack organization, unimanual scooping and bimanual pouring—under both ID and OOD settings (Figure 5), using the LocoMan robot and data collected via the Human2LocoMan system. Success rates and task scores are used as evaluation metrics, with HIT [20] and HPT [32] serving as baselines. For detailed information on the experimental setups, including data statistics, model hyperparameters, masking strategies for embodiment alignment, and training configurations for MXT and the baselines, please refer to Appendix Section IV-F.

# B. Results and Analysis

(1) Does the Human2LocoMan system enable versatile quadrupedal manipulation capabilities? Please refer to Appendix Section IV-G for detailed analysis.

TABLE I: Result Summary. We report success rate (SR) $\uparrow$ in % and task score (TS) $\uparrow$ for each task.	We highlight the best
performance in <b>bold</b> and the second best in <u>underline</u> .	

				Toy Collection				Shoe Rack Organization				Scooping			Pouring											
				Unin	nanual			Bim	anual			Unim	anual			Bim	anual			Unim	anual			Bim	nual	
			II	)	OC	D	II	)	00	DD	П	D	OC	D	II	)	OC	D	П	D	00	D	II	)	OC	D
Method	Pretrained	Data	SR	TS	SR	TS	SR	TS	SR	TS	SR	TS	SR	TS	SR	TS	SR	TS	SR	TS	SR	TS	SR	TS	SR	TS
HIT HIT	-	smaller larger	54.2 79.2	42 57	41.6 58.3	15 23	45.8 58.3	37 47	41.6 58.3	16 21	<u>87.5</u> 79.2	$\tfrac{112}{107}$	75.0 <b>83.3</b>	50 52	66.7 83.3	52 63	25.0 33.3	14 15	58.3 66.7	96 106	16.7 33.3	30 34	58.3 70.8	62 72	16.7 8.33	17 7
MXT	Ν	smaller	70.8	56	33.3	20	66.7	54	41.7	15	87.5	109	16.7	10	66.7	52	33.3	14	62.5	105	16.7	30	75.0	75	33.3	24
MXT	N	larger	87.5	67	83.3	31	70.8	53	41.7	16	83.3	107	50.0	37	75.0	60	58.3	23	62.5	98	41.7	38	79.2	76	33.3	22
MXT	Y	smaller	91.7	66	83.3	30	83.3	62	83.3	31	83.3	103	75.0	47	79.2	61	58.3	24	87.5	129	25.0	35	83.3	83	58.3	<u>33</u>
MXT	Y	larger	95.8	67	91.7	34	91.7	67	100	36	95.8	116	83.3	52	83.3	63	75.0	29	87.5	129	66.7	52	91.7	88	83.3	42

\* Number of trajectories: TC-Uni smaller=20, larger=40; TC-Bi smaller=30, larger=60; SO-Uni smaller=40, larger=80; SO-Bi smaller=40, larger=80; Scoop-Uni smaller=30, larger=60; Pour-Bi smaller=30, larger=60.

(2) How does MXT compare to state-of-the-art imitation learning architectures? Compared to HIT. As shown in Table I, MXT without pretraining achieves comparable or better performance than HIT across unimanual and bimanual tasks, under both ID and OOD settings. With pretraining, MXT consistently outperforms HIT in both success rate and task score. From Figure 7, MXT shows lower validation loss on most tasks, indicating better convergence and scalability with larger datasets. Although HIT achieves lower validation loss in some cases, it performs comparably only in tasks like shoe organization, where limited object variation may favor HIT despite its lack of modularity (Figure 6). Compared to HPT. As shown in Figures 4 and 8, MXT consistently surpasses HPT across all pretraining and data size settings on toy collection tasks. MXT also exhibits better validation loss trends and avoids the severe overfitting observed in HPT, highlighting the generalization benefits of its modular architecture.

(3) How does human data collected by Human2LocoMan contribute to imitation learning performance? Efficiency, robustness, and generalizability. As shown in Table I, pretraining on human data has a substantial positive impact on LocoMan manipulation performance. The policy maintains strong performance even when robot data is limited, highlighting both its efficiency and robustness. Specifically, comparing MXT-Pretrained to MXT-Scratch in Table I, we observe that pretraining improves performance on TC-Uni, TC-Bi, and Scooping tasks under ID settings, where objects exhibit diverse locations. MXT-Pretrained tends to produce smoother and more robust motions, enabling more accurate localization of target objects. For instance, as shown in Figure 6, MXT-Pretrained achieves substantially better scooping performance-which requires precise localization-compared to all other methods. Moreover, Table I reveals large performance gaps on OOD objects in tasks such as TC-Bi, SO-Uni, and Pouring, where OOD objects differ significantly from ID objects in shape, texture, and color. These results suggest that MXT, by leveraging human demonstrations during the pretraining stage, is able to generalize effectively to novel scenarios unseen during robot training.

(4) Do the design choices in MXT facilitate positive transfer from Human to LocoMan? Our framework presents positive cross-embodiment transfer despite substantial embodiment gaps. For comparisons, we provide SR and TS results from 36 trials in Fig. 4(b). HPT performs consistently worse than



Fig. 4: Ablation study on unimanual and bimanual toy collection. We compare MXT, its ablation MXT-Agg, and baseline HPT on SR and TS. Here, 100 denotes the larger training set (40 trajectories for TC-Uni, 60 for TC-Bi), while 50 denotes the smaller set (20 for TC-Uni, 30 for TC-Bi).

MXT, both when finetuned and trained from scratch. We attribute part of this performance gap to HPT using frozen image encoders by default. We also provide additional ablations of MXT where we *aggregate* the input modalities, tokenize them with a single tokenizer, and decode actions with a single detokenizer; this baseline (marked with "Agg" in Fig. 4(b)) incorporates the key HPT designs, while finetuning the vision encoders and remaining architecturally comparable to MXT. MXT consistently benefits from pretraining and outperforms this baseline when both are finetuned, highlighting the advantage of modularized tokenization for leveraging human data. Notably, MXT-Agg sometimes transfers suboptimally, likely due to the lack of modular design and the trade-off between improved performance from adaptive tokenization (image encoder finetuning) and reduced trunk transferability.

# IV. CONCLUSIONS

In this paper, we introduce Human2LocoMan, a unified framework for flexible data collection and cross-embodiment learning to enable versatile quadrupedal manipulation skills on the open-source LocoMan platform. Our teleoperation and human data collection systems allow efficient acquisition of large-scale, high-quality datasets by bridging the action spaces between human and robot embodiments. Built upon this foundation, we propose Modularized Cross-embodiment Transformer, a modular policy architecture that supports positive transfer from human demonstrations to robot policies. Through extensive experiments on six challenging household tasks, we demonstrate that Human2LocoMan enables strong performance, efficient training, and robust generalization to out-of-distribution scenarios, outperforming strong imitation learning baselines. Our results highlight the effectiveness of cross-embodiment learning and modularized policy design in advancing scalable, versatile quadrupedal manipulation.

#### REFERENCES

- F. Jenelten, J. He, F. Farshidian, and M. Hutter, "Dtc: Deep tracking control," *Science Robotics*, vol. 9, no. 86, p. eadh5401, 2024.
- [2] S. Choi, G. Ji, J. Park, H. Kim, J. Mun, J. H. Lee, and J. Hwangbo, "Learning quadrupedal locomotion on deformable terrain," *Science Robotics*, vol. 8, no. 74, p. eade2256, 2023.
- [3] J. Lee, J. Hwangbo, L. Wellhausen, V. Koltun, and M. Hutter, "Learning quadrupedal locomotion over challenging terrain," *Science robotics*, vol. 5, no. 47, p. eabc5986, 2020.
- [4] R. Yang, G. Yang, and X. Wang, "Neural volumetric memory for visual locomotion control," in CVPR 2023, 2023.
- [5] Y. Yang, G. Shi, C. Lin, X. Meng, R. Scalise, M. G. Castro, W. Yu, T. Zhang, D. Zhao, J. Tan *et al.*, "Agile continuous jumping in discontinuous terrains," *arXiv preprint arXiv:2409.10923*, 2024.
- [6] A. Kumar, Z. Fu, D. Pathak, and J. Malik, "Rma: Rapid motor adaptation for legged robots," arXiv preprint arXiv:2107.04034, 2021.
- [7] B. Lindqvist, S. Karlsson, A. Koval, I. Tevetzidis, J. Haluška, C. Kanellakis, A.-a. Agha-mohammadi, and G. Nikolakopoulos, "Multimodality robotic systems: Integrated combined legged-aerial mobility for subterranean search-and-rescue," *Robotics and Autonomous Systems*, 2022.
- [8] Z. Fu, X. Cheng, and D. Pathak, "Deep whole-body control: learning a unified policy for manipulation and locomotion," in *Conference on Robot Learning*. PMLR, 2023, pp. 138–149.
- [9] Q. Wu, Z. Fu, X. Cheng, X. Wang, and C. Finn, "Helpful doggybot: Open-world object fetching using legged robots and vision-language models," arXiv preprint arXiv:2410.00231, 2024.
- [10] H. Ha, Y. Gao, Z. Fu, J. Tan, and S. Song, "Umi on legs: Making manipulation policies mobile with manipulation-centric whole-body controllers," arXiv preprint arXiv:2407.10353, 2024.
- [11] M. Seo, H. A. Park, S. Yuan, Y. Zhu, and L. Sentis, "Legato: Cross-embodiment imitation using a grasping tool," *arXiv preprint* arXiv:2411.03682, 2024.
- [12] X. He, C. Yuan, W. Zhou, R. Yang, D. Held, and X. Wang, "Visual manipulation with legs," arXiv preprint arXiv:2410.11345, 2024.
- [13] R.-Z. Qiu, Y. Song, X. Peng, S. A. Suryadevara, G. Yang, M. Liu, M. Ji, C. Jia, R. Yang, X. Zou *et al.*, "Wildlma: Long horizon locomanipulation in the wild," *arXiv preprint arXiv:2411.15131*, 2024.
- [14] C. Lin, X. Liu, Y. Yang, Y. Niu, W. Yu, T. Zhang, J. Tan, B. Boots, and D. Zhao, "Locoman: Advancing versatile quadrupedal dexterity with lightweight loco-manipulators," arXiv preprint arXiv:2403.18197, 2024.
- [15] S. Schaal, "Learning from demonstration," Advances in neural information processing systems, vol. 9, 1996.
- [16] Y. Zhu, A. Joshi, P. Stone, and Y. Zhu, "Viola: Imitation learning for vision-based manipulation with object proposal priors," in *Conference* on Robot Learning. PMLR, 2023, pp. 1199–1210.
- [17] C. Wang, L. Fan, J. Sun, R. Zhang, L. Fei-Fei, D. Xu, Y. Zhu, and A. Anandkumar, "Mimicplay: Long-horizon imitation learning by watching human play," *arXiv preprint arXiv:2302.12422*, 2023.
- [18] T. Z. Zhao, V. Kumar, S. Levine, and C. Finn, "Learning fine-grained bimanual manipulation with low-cost hardware," *arXiv preprint* arXiv:2304.13705, 2023.
- [19] P. Wu, Y. Shentu, Z. Yi, X. Lin, and P. Abbeel, "Gello: A general, low-cost, and intuitive teleoperation framework for robot manipulators," in 2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). IEEE, 2024, pp. 12156–12163.
- [20] Z. Fu, Q. Zhao, Q. Wu, G. Wetzstein, and C. Finn, "Humanplus: Humanoid shadowing and imitation from humans," arXiv preprint arXiv:2406.10454, 2024.
- [21] X. Cheng, J. Li, S. Yang, G. Yang, and X. Wang, "Open-television: Teleoperation with immersive active visual feedback," *arXiv preprint* arXiv:2407.01512, 2024.
- [22] T. He, Z. Luo, X. He, W. Xiao, C. Zhang, W. Zhang, K. Kitani, C. Liu, and G. Shi, "Omnih2o: Universal and dexterous humanto-humanoid whole-body teleoperation and learning," *arXiv preprint arXiv:2406.08858*, 2024.
- [23] S. Yang, M. Liu, Y. Qin, R. Ding, J. Li, X. Cheng, R. Yang, S. Yi, and X. Wang, "Ace: A cross-platform visual-exoskeletons system for low-cost dexterous teleoperation," *arXiv preprint arXiv:2408.11805*, 2024.
- [24] K. Bousmalis, A. Irpan, P. Wohlhart, Y. Bai, M. Kelcey, M. Kalakrishnan, L. Downs, J. Ibarz, P. Pastor, K. Konolige *et al.*, "Using simulation and domain adaptation to improve efficiency of deep robotic grasping,"

in 2018 IEEE international conference on robotics and automation (ICRA). IEEE, 2018, pp. 4243–4250.

- [25] Z. Jiang, Y. Xie, K. Lin, Z. Xu, W. Wan, A. Mandlekar, L. Fan, and Y. Zhu, "Dexmimicgen: Automated data generation for bimanual dexterous manipulation via imitation learning," *arXiv preprint* arXiv:2410.24185, 2024.
- [26] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, X. Chen, K. Choromanski, T. Ding, D. Driess, A. Dubey, C. Finn *et al.*, "Rt-2: Visionlanguage-action models transfer web knowledge to robotic control," *arXiv preprint arXiv:2307.15818*, 2023.
- [27] C. Wang, H. Shi, W. Wang, R. Zhang, L. Fei-Fei, and C. K. Liu, "Dexcap: Scalable and portable mocap data collection system for dexterous manipulation," arXiv preprint arXiv:2403.07788, 2024.
- [28] H. Bharadhwaj, A. Gupta, V. Kumar, and S. Tulsiani, "Towards generalizable zero-shot manipulation via translating human interaction plans," in 2024 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2024, pp. 6904–6911.
- [29] M. K. Srirama, S. Dasari, S. Bahl, and A. Gupta, "Hrp: Human affordances for robotic pre-training," *arXiv preprint arXiv:2407.18911*, 2024.
- [30] S. Chen, C. Wang, K. Nguyen, L. Fei-Fei, and C. K. Liu, "Arcap: Collecting high-quality human demonstrations for robot learning with augmented reality feedback," arXiv preprint arXiv:2410.08464, 2024.
- [31] S. Kareer, D. Patel, R. Punamiya, P. Mathur, S. Cheng, C. Wang, J. Hoffman, and D. Xu, "Egomimic: Scaling imitation learning via egocentric video," arXiv preprint arXiv:2410.24221, 2024.
- [32] L. Wang, X. Chen, J. Zhao, and K. He, "Scaling proprioceptive-visual learning with heterogeneous pre-trained transformers," 2024. [Online]. Available: https://arxiv.org/abs/2409.20537
- [33] J. Carpentier, G. Saurel, G. Buondonno, J. Mirabel, F. Lamiraux, O. Stasse, and N. Mansard, "The pinocchio c++ library: A fast and flexible implementation of rigid body dynamics algorithms and their analytical derivatives," in 2019 IEEE/SICE International Symposium on System Integration (SII). IEEE, 2019, pp. 614–619.
- [34] L. Wang, X. Chen, J. Zhao, and K. He, "Scaling proprioceptive-visual learning with heterogeneous pre-trained transformers," *arXiv preprint arXiv:2409.20537*, 2024.
- [35] Y. Ji, G. B. Margolis, and P. Agrawal, "Dribblebot: Dynamic legged manipulation in the wild," in 2023 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2023, pp. 5155–5162.
- [36] F. Shi, T. Homberger, J. Lee, T. Miki, M. Zhao, F. Farshidian, K. Okada, M. Inaba, and M. Hutter, "Circus anymal: A quadruped learning dexterous manipulation with its limbs," in 2021 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2021, pp. 2316–2323.
- [37] X. Cheng, A. Kumar, and D. Pathak, "Legs as manipulator: Pushing quadrupedal agility beyond locomotion," in 2023 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2023, pp. 5106–5112.
- [38] M. Zhang, Y. Ma, T. Miki, and M. Hutter, "Learning to open and traverse doors with a legged manipulator," arXiv preprint arXiv:2409.04882, 2024.
- [39] Z. He, K. Lei, Y. Ze, K. Sreenath, Z. Li, and H. Xu, "Learning visual quadrupedal loco-manipulation from demonstrations," *arXiv preprint* arXiv:2403.20328, 2024.
- [40] T. Huang, N. Sontakke, K. N. Kumar, I. Essa, S. Nikolaidis, D. W. Hong, and S. Ha, "Bayrntune: Adaptive bayesian domain randomization via strategic fine-tuning," *arXiv preprint arXiv:2310.10606*, 2023.
- [41] S. Jeon, M. Jung, S. Choi, B. Kim, and J. Hwangbo, "Learning whole-body manipulation for quadrupedal robot," *IEEE Robotics and Automation Letters*, vol. 9, no. 1, pp. 699–706, 2023.
- [42] J. Stolle, P. Arm, M. Mittal, and M. Hutter, "Perceptive pedipulation with local obstacle avoidance," arXiv preprint arXiv:2409.07195, 2024.
- [43] R.-Z. Qiu, Y. Hu, G. Yang, Y. Song, Y. Fu, J. Ye, J. Mu, R. Yang, N. Atanasov, S. Scherer *et al.*, "Learning generalizable feature fields for mobile manipulation," *arXiv preprint arXiv:2403.07563*, 2024.
- [44] K. N. Kumar, I. Essa, and S. Ha, "Cascaded compositional residual learning for complex interactive behaviors," *IEEE Robotics and Automation Letters*, vol. 8, no. 8, pp. 4601–4608, 2023.
- [45] Y. Feng, C. Hong, Y. Niu, S. Liu, Y. Yang, W. Yu, T. Zhang, J. Tan, and D. Zhao, "Learning multi-agent loco-manipulation for long-horizon quadrupedal pushing," arXiv preprint arXiv:2411.07104, 2024.
- [46] Z. Xiong, B. Chen, S. Huang, W.-W. Tu, Z. He, and Y. Gao,

"Mqe: Unleashing the power of interaction with multi-agent quadruped environment," *arXiv preprint arXiv:2403.16015*, 2024.

- [47] T. An, J. Lee, M. Bjelonic, F. De Vincenti, and M. Hutter, "Solving multi-entity robotic problems using permutation invariant neural networks," *arXiv preprint arXiv:2402.18345*, 2024.
- [48] O. Nachum, M. Ahn, H. Ponte, S. Gu, and V. Kumar, "Multiagent manipulation via locomotion using hierarchical sim2real," *arXiv* preprint arXiv:1908.05224, 2019.
- [49] Y. Ji, B. Zhang, and K. Sreenath, "Reinforcement learning for collaborative quadrupedal manipulation of a payload over challenging terrain," in 2021 IEEE 17th International Conference on Automation Science and Engineering (CASE). IEEE, 2021, pp. 899–904.
- [50] G. Pan, Q. Ben, Z. Yuan, G. Jiang, Y. Ji, J. Pang, H. Liu, and H. Xu, "Roboduet: A framework affording mobile-manipulation and crossembodiment," *arXiv preprint arXiv:2403.17367*, 2024.
- [51] M. Liu, Z. Chen, X. Cheng, Y. Ji, R. Yang, and X. Wang, "Visual whole-body control for legged loco-manipulation," *arXiv preprint* arXiv:2402.16796, 2024.
- [52] J. Zhang, N. Gireesh, J. Wang, X. Fang, C. Xu, W. Chen, L. Dai, and H. Wang, "Gamma: Graspability-aware mobile manipulation policy learning based on online grasping pose fusion," in 2024 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2024, pp. 1399–1405.
- [53] P. Arm, M. Mittal, H. Kolvenbach, and M. Hutter, "Pedipulate: Enabling manipulation skills using a quadruped robot's leg," in 41st IEEE Conference on Robotics and Automation (ICRA 2024), 2024.
- [54] B. D. Argall, S. Chernova, M. Veloso, and B. Browning, "A survey of robot learning from demonstration," *Robotics and autonomous* systems, vol. 57, no. 5, pp. 469–483, 2009.
- [55] A. Hussein, M. M. Gaber, E. Elyan, and C. Jayne, "Imitation learning: A survey of learning methods," ACM Computing Surveys (CSUR), vol. 50, no. 2, pp. 1–35, 2017.
- [56] H. Ravichandar, A. S. Polydoros, S. Chernova, and A. Billard, "Recent advances in robot learning from demonstration," *Annual review of control, robotics, and autonomous systems*, vol. 3, no. 1, pp. 297– 330, 2020.
- [57] M. Memmel, J. Berg, B. Chen, A. Gupta, and J. Francis, "Strap: Robot sub-trajectory retrieval for augmented policy learning," in *The Thirteenth International Conference on Learning Representations*, 2025.
- [58] A. Khazatsky, K. Pertsch, S. Nair, A. Balakrishna, S. Dasari, S. Karamcheti, S. Nasiriany, M. K. Srirama, L. Y. Chen, K. Ellis *et al.*, "Droid: A large-scale in-thewild robot manipulation dataset," *arXiv preprint arXiv:2403.12945*, 2024.
- [59] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Hausman, A. Herzog, J. Hsu *et al.*, "Rt-1: Robotics transformer for real-world control at scale," *arXiv preprint arXiv:2212.06817*, 2022.
- [60] N. M. M. Shafiullah, A. Rai, H. Etukuru, Y. Liu, I. Misra, S. Chintala, and L. Pinto, "On bringing robots home," arXiv preprint arXiv:2311.16098, 2023.
- [61] X. Lin, J. So, S. Mahalingam, F. Liu, and P. Abbeel, "Spawnnet: Learning generalizable visuomotor skills from pre-trained network," in 2024 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2024, pp. 4781–4787.
- [62] H. Liu, S. Nasiriany, L. Zhang, Z. Bao, and Y. Zhu, "Robot learning on the job: Human-in-the-loop autonomy and learning during deployment," *The International Journal of Robotics Research*, p. 02783649241273901, 2022.
- [63] T. He, Z. Luo, W. Xiao, C. Zhang, K. Kitani, C. Liu, and G. Shi, "Learning human-to-humanoid real-time whole-body teleoperation," arXiv preprint arXiv:2403.04436, 2024.
- [64] Z. Zhang, Y. Niu, Z. Yan, and S. Lin, "Real-time whole-body imitation by humanoid robots and task-oriented teleoperation using an analytical mapping method and quantitative evaluation," *Applied Sciences*, vol. 8, no. 10, p. 2005, 2018.
- [65] L. P. Poubel, S. Sakka, D. Ćehajić, and D. Creusot, "Support changes during online human motion imitation by a humanoid robot using task specification," in 2014 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2014, pp. 1782–1787.
- [66] B. Sen, M. Wang, N. Thakur, A. Agarwal, and P. Agrawal, "Learning to look around: Enhancing teleoperation and learning with a humanlike actuated neck," *arXiv preprint arXiv:2411.00704*, 2024.
- [67] C. Lu, X. Cheng, J. Li, S. Yang, M. Ji, C. Yuan, G. Yang, S. Yi, and

X. Wang, "Mobile-television: Predictive motion priors for humanoid whole-body control," *arXiv preprint arXiv:2412.07773*, 2024.

- [68] R. Ding, Y. Qin, J. Zhu, C. Jia, S. Yang, R. Yang, X. Qi, and X. Wang, "Bunny-visionpro: Real-time bimanual dexterous teleoperation for imitation learning," arXiv preprint arXiv:2407.03162, 2024.
- [69] H. Fang, H.-S. Fang, Y. Wang, J. Ren, J. Chen, R. Zhang, W. Wang, and C. Lu, "Airexo: Low-cost exoskeletons for learning whole-arm manipulation in the wild," in 2024 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2024, pp. 15031–15038.
- [70] C. Chi, Z. Xu, C. Pan, E. Cousineau, B. Burchfiel, S. Feng, R. Tedrake, and S. Song, "Universal manipulation interface: In-the-wild robot teaching without in-the-wild robots," *arXiv preprint arXiv:2402.10329*, 2024.
- [71] G. Papagiannis, N. Di Palo, P. Vitiello, and E. Johns, "R+ x: Retrieval and execution from everyday human videos," arXiv preprint arXiv:2407.12957, 2024.
- [72] M. J. Kim, K. Pertsch, S. Karamcheti, T. Xiao, A. Balakrishna, S. Nair, R. Rafailov, E. Foster, G. Lam, P. Sanketi *et al.*, "Openvla: An open-source vision-language-action model," *arXiv preprint arXiv:2406.09246*, 2024.
- [73] O. M. Team, D. Ghosh, H. Walke, K. Pertsch, K. Black, O. Mees, S. Dasari, J. Hejna, T. Kreiman, C. Xu *et al.*, "Octo: An open-source generalist robot policy," *arXiv preprint arXiv:2405.12213*, 2024.
- [74] A. O'Neill, A. Rehman, A. Maddukuri, A. Gupta, A. Padalkar, A. Lee, A. Pooley, A. Gupta, A. Mandlekar, A. Jain *et al.*, "Open x-embodiment: Robotic learning datasets and rt-x models: Open xembodiment collaboration 0," in 2024 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2024, pp. 6892–6903.
- [75] L. Y. Chen, K. Hari, K. Dharmarajan, C. Xu, Q. Vuong, and K. Goldberg, "Mirage: Cross-embodiment zero-shot policy transfer with crosspainting," arXiv preprint arXiv:2402.19249, 2024.
- [76] R. Doshi, H. Walke, O. Mees, S. Dasari, and S. Levine, "Scaling cross-embodied learning: One policy for manipulation, navigation, locomotion and aviation," arXiv preprint arXiv:2408.11812, 2024.

#### APPENDIX

# A. Human2LocoMan Motion Mapping

We map the human wrist motions to LocoMan's endeffector motions, map the human head motions to LocoMan's torso motions, and hand poses to LocoMan's gripper actions. The 6D poses of the human hand, head, and wrist poses in SE(3) in the VR-defined world frame are streamed from the VR set to the Human2LocoMan teleoperation server. The human head pose is represented as  $(\boldsymbol{x}_{vr}^{head}, \boldsymbol{R}_{vr}^{head})$ , and the wrist poses are  $(\boldsymbol{x}_{vr}^{r.wrist}, \boldsymbol{R}_{vr}^{r.wrist})$  and  $(\boldsymbol{x}_{vr}^{l.wrist}, \boldsymbol{R}_{vr}^{l.wrist})$ , where  $\boldsymbol{x}_{vr}$  denotes the translation and  $\boldsymbol{R}_{vr}$  denotes the rotation in the VR-defined world frame. Then, the 6D poses can be transformed into the unified frame  $\mathcal{F}_u$   $(\boldsymbol{x}_{uni}, \boldsymbol{R}_{uni}) =$  $(\boldsymbol{R}_{uni}^{vr}, \boldsymbol{R}_{uni}^{vr}, \boldsymbol{R}_{vr}^{un}, \boldsymbol{R}_{vr})$ , where  $\boldsymbol{R}_{uni}^{vr}$  is the rotation matrix of the VR-defined frame relative to the unified frame  $\mathcal{F}_u$ .

To initialize the teleoperation for each manipulation mode, the robot is transferred to a reset pose randomly initialized within a small range, termed as  $p_0 = (x_{\text{uni},0}^{\text{troso}}, R_{\text{uni},0}^{\text{treef}}, R_{\text{uni},0}^{\text{treef}}, R_{\text{uni},0}^{\text{leef}}, R_{\text{uni},0}^{\text{leef}}, \theta_0^{\text{gripper}})$ , including the 6D poses of the torso and both end effectors, and the gripper angles. The human operator starts to teleoperate the robot after a initializing posture. The target pose for the robot at time step t,  $p_t^t = (x_{\text{uni},t}^{\text{torso,t}}, R_{\text{uni},t}^{\text{torso,t}}, x_{\text{uni},t}^{\text{reef,t}}, R_{\text{uni},t}^{\text{l-eef,t}}, R_{\text{uni},t}^{\text{l-eef,t}}, \theta_t^{\text{gripper},t})$ , can be expressed as follows.

$$\begin{aligned} \boldsymbol{x}_{\text{uni},t}^{\text{torso,t}} &= \boldsymbol{x}_{\text{uni},0}^{\text{torso}} + \alpha^{\text{torso}}(\boldsymbol{x}_{\text{uni},t}^{\text{head}} - \boldsymbol{x}_{\text{uni},0}^{\text{head}}) \\ \boldsymbol{R}_{\text{uni},t}^{\text{torso,t}} &= \boldsymbol{R}_{\text{uni},0}^{\text{torso}}((\boldsymbol{R}_{\text{uni},0}^{\text{head}})^{\top}\boldsymbol{R}_{\text{uni},t}^{\text{head}}) \\ \boldsymbol{x}_{\text{uni},t}^{\text{reef,t}} &= \boldsymbol{x}_{\text{uni},0}^{\text{reef}} + \alpha^{\text{reef}}(\boldsymbol{x}_{\text{uni},t}^{\text{r-wrist}} - \boldsymbol{x}_{\text{uni},0}^{\text{torso,t}}) \\ \boldsymbol{R}_{\text{uni},t}^{\text{reef,t}} &= \boldsymbol{R}_{\text{uni},0}^{\text{reef}}((\boldsymbol{R}_{\text{uni},0}^{\text{r-wrist}})^{\top}\boldsymbol{R}_{\text{uni},t}^{\text{r-wrist}}) \\ \boldsymbol{x}_{\text{uni},t}^{\text{l-eef,t}} &= \boldsymbol{x}_{\text{uni},0}^{\text{l-eef}} + \alpha^{\text{l-eef}}(\boldsymbol{x}_{\text{uni},t}^{\text{l-wrist}} - \boldsymbol{x}_{\text{uni},0}^{\text{l-wrist}}) \\ \boldsymbol{R}_{\text{uni},t}^{\text{l-eef,t}} &= \boldsymbol{R}_{\text{uni},0}^{\text{l-eef}}((\boldsymbol{R}_{\text{uni},0}^{\text{l-wrist}})^{\top}\boldsymbol{R}_{\text{uni},t}^{\text{l-wrist}}) \\ \boldsymbol{\theta}_{t}^{\text{gripper,t}} &= \frac{\boldsymbol{\theta}_{\text{max}}^{\text{gripper}} - \boldsymbol{\theta}_{\text{min}}^{\text{gripper}}}{\boldsymbol{d}_{\text{max}}^{\text{torg}}} \circ \boldsymbol{d}_{t}^{\text{tip}} + \boldsymbol{\theta}_{\text{min}}^{\text{gripper}} \end{aligned}$$

Here,  $\alpha^{\text{torso}}$ ,  $\alpha^{\text{r-eef}}$ , and  $\alpha^{\text{l-eef}}$ , are the scaling factors to map human's motions to robot's torso, right end effector, and left end effector, respectively.  $x_{\text{max}}^{\text{gripper}}$  and  $x_{\text{min}}^{\text{gripper}}$  are the maximum and minimum gripper angles, respectively.  $d_t^{\text{tip}}$ represents the distances between the reference finger tips of both human hands at time step t, and  $d_{\text{max}}^{\text{tip}}$  is the maximum finger tip distance for the human operator.

#### B. Human2LocoMan Whole-Body Controller

The robot target pose at time t,  $p_t^t$ , is calculated from the teleoperation server, and sent to the whole-body controller of the LocoMan robot, which is adapted from the one introduced in [14], a unified whole-body controller designed to track the desired poses of the torso, end effectors, and feet across multiple operation modes. We employ null-space projection for kinematic tracking and quadratic programming for dynamic optimization to compute the desired joint positions, velocities, and torques. To handle the large embodiment gap between the human and the LocoMan robots, and to facilitate smooth teleoperation of a dynamic quadrupedal

platform with whole-body motions, we consider the handling and recovery from robot's joint limits, singularity, and selfcollision, and fast motions. We compute the manipulability index as:

$$I_{\text{mani}} = \sqrt{\det(\mathbf{J}\mathbf{J}^{\top})} \tag{2}$$

to assess the proximity of the target pose to singularity, where **J** represents the Jacobian of the robot's manipulator at the target pose. If  $I_{\text{mani}}$  falls below a predefined threshold  $\tau_{\text{mani}}$ , the target pose is considered near singularity. To detect self-collisions, we utilize the Pinocchio library [33] to compute collision pairs among the robot's body parts. If any of the following conditions are met—joint limit violation, singularity, or self-collision—the whole-body controller tracks  $p_{t-1}^{t}$  instead of  $p_t^{t}$ . To mitigate rapid movements, we apply linear interpolation between  $x_{\text{uni},t}^{\text{torso,t}}$  and  $\theta_{t-1}^{\text{reef,t}}$  and  $x_{\text{uni},t-1}^{\text{reef,t}}$ , and  $x_{\text{uni},t-1}^{\text{treef,t}}$ , and  $R_{\text{uni},t-1}^{\text{treef,t}}$ , and  $R_{\text{uni},t-$ 

## C. Human2LocoMan Data Collection

We record the robot data  $\{\mathcal{D}_t^R\}_{t=1}^T$  during teleoperation, where  $\mathcal{D}_t^R = \{o_t^R, a_t^R\}$  is the robot data at time step tincluding the robot observations  $o_t^R$  and robot actions  $a_t^R$ , and T is the episode length. We define the  $I_{\min,t}^R$  and  $I_{\text{wrist},t}^R$  are images obtained from the robot's main stereo camera and the wrist camera, respectively. Then, we can formulate  $o_t^R$  and  $a_t^R$  in the dataset as follows.

$$\boldsymbol{o}_{t}^{\mathrm{R}}[\text{main image}] := I_{\text{main},t},$$
  

$$\boldsymbol{o}_{t}^{\mathrm{R}}[\text{wrist image}] := I_{\text{wrist},t},$$
  

$$\boldsymbol{o}_{t}^{\mathrm{R}}[\text{body pose}] := [\boldsymbol{x}_{\text{uni},t}^{\text{torso}}, \boldsymbol{R}_{\text{uni},t}^{\text{torso}}],$$
  

$$\boldsymbol{o}_{t}^{\mathrm{R}}[\text{EEF pose}] := [\boldsymbol{x}_{\text{uni},t}^{\text{r-eef}}, \boldsymbol{R}_{\text{uni},t}^{\text{t-eef}}, \boldsymbol{x}_{\text{uni},t}^{\text{l-eef}}],$$
  

$$\boldsymbol{o}_{t}^{\mathrm{R}}[\text{EEF to body pose}] := [\boldsymbol{x}_{\text{uni},t}^{\text{r-eef}} - \boldsymbol{x}_{\text{uni},t}^{\text{torso}}, (\boldsymbol{R}_{\text{uni},t}^{\text{torso}})^{\top} \boldsymbol{R}_{\text{uni},t}^{\text{t-eef}}],$$
  

$$\boldsymbol{v}_{t}^{\mathrm{R}}[\text{gripper angles}] := \boldsymbol{\theta}_{t}^{\text{gripper}},$$
  

$$\boldsymbol{a}_{t}^{\mathrm{R}}[\text{body pose}] := [\boldsymbol{x}_{\text{uni},t}^{\text{torso}, t}, \boldsymbol{R}_{\text{uni},t}^{\text{torso}, t}],$$
  

$$\boldsymbol{a}_{t}^{\mathrm{R}}[\text{EEF pose}] := [\boldsymbol{x}_{\text{uni},t}^{\text{r-eef}, t}, \boldsymbol{x}_{\text{uni},t}^{\text{t-eef}, t}, \boldsymbol{R}_{\text{uni},t}^{\text{t-eef}, t}],$$
  

$$\boldsymbol{a}_{t}^{\mathrm{R}}[\text{gripper angles}] := \boldsymbol{\theta}_{t}^{\text{gripper}, t}$$
  
(3)

We record the human data  $\{\mathcal{D}_t^H\}_{t=1}^T$  in real time during human's manipulation. Similarly, the human data at time step  $t \mathcal{D}_t^H = \{o_t^H, a_t^H\}$  can be defined by human observations  $o_t^H$ 

and human actions  $a_t^{\rm H}$  as follows.

$$\boldsymbol{o}_{t}^{\mathrm{H}}[\text{main image}] := I_{\mathrm{main},t}^{\mathrm{H}},$$
  

$$\boldsymbol{o}_{t}^{\mathrm{H}}[\text{body pose}] := [\boldsymbol{x}_{\mathrm{uni},t}^{\mathrm{head}}, \boldsymbol{R}_{\mathrm{uni},t}^{\mathrm{head}}],$$
  

$$\boldsymbol{o}_{t}^{\mathrm{H}}[\text{EEF pose}] := [\boldsymbol{x}_{\mathrm{uni},t}^{\mathrm{r-wrist}}, \boldsymbol{R}_{\mathrm{uni},t}^{\mathrm{l-wrist}}, \boldsymbol{x}_{\mathrm{uni},t}^{\mathrm{l-wrist}}, \boldsymbol{R}_{\mathrm{uni},t}^{\mathrm{l-wrist}}],$$
  

$$\boldsymbol{o}_{t}^{\mathrm{H}}[\text{EEF to body pose}] := [\boldsymbol{x}_{\mathrm{uni},t}^{\mathrm{r-wrist}} - \boldsymbol{x}_{\mathrm{uni},t}^{\mathrm{head}}, (\boldsymbol{R}_{\mathrm{uni},t}^{\mathrm{head}})^{\top} \boldsymbol{R}_{\mathrm{uni},t}^{\mathrm{r-wrist}}],$$
  

$$\boldsymbol{o}_{t}^{\mathrm{H}}[\text{EEF to body pose}] := [\boldsymbol{x}_{\mathrm{uni},t}^{\mathrm{r-wrist}} - \boldsymbol{x}_{\mathrm{uni},t}^{\mathrm{head}}, (\boldsymbol{R}_{\mathrm{uni},t}^{\mathrm{head}})^{\top} \boldsymbol{R}_{\mathrm{uni},t}^{\mathrm{l-wrist}}],$$
  

$$\boldsymbol{v}_{t}^{\mathrm{H}}[\text{grasping states}] := \boldsymbol{\theta}_{t}^{\mathrm{gripper}},$$
  

$$\boldsymbol{a}_{t}^{\mathrm{H}}[\text{body pose}] := [\boldsymbol{x}_{\mathrm{uni},t}^{\mathrm{head}, t}, \boldsymbol{R}_{\mathrm{uni},t}^{\mathrm{r-wrist}, t}],$$
  

$$\boldsymbol{a}_{t}^{\mathrm{H}}[\text{EEF pose}] := [\boldsymbol{x}_{\mathrm{uni},t}^{\mathrm{r-wrist}, t}, \boldsymbol{R}_{\mathrm{uni},t}^{\mathrm{r-wrist}, t}],$$
  

$$\boldsymbol{a}_{t}^{\mathrm{H}}[\text{grasping actions}] := \boldsymbol{\theta}_{t}^{\mathrm{gripper}, t}$$
  
(4)

#### D. Design Details of MXT

**Tokenizers.** The tokenizers T transform raw observations into tokens for the transformer trunk. Drawing from the design in previous works [34], we use a cross attention layer to format observational features into a fixed number of tokens. For image inputs, the features are obtained from a pretrained ResNet encoder that can be finetuned during training; for proprioceptive or state-like inputs, the features are computed by passing the raw input through a trainable MLP network.

**Detokenizers.** The detokenizers D serve as action decoder heads and map output tokens from the trunk to actions in each embodiment's action space. We adopt the action chunking technique [18]. At each inference step, the detokenizers predict an action sequence of h steps and temporal ensemble is applied to the outputs, following [18]. Within each detokenizer, we use a cross attention layer to transform the latent action tokens output by the trunk to a sequence of actions with length h and appropriate action dimensions.

**Trunk.** The trunk is an encoder-decoder transformer, where the input sequence length and the output sequence length are both fixed, as the number of tokens for each input or output modality is fixed by design. By sharing the trunk weights across the human and robot embodiments, the trunk is trained to capture the common decision making patterns across different embodiments.

**Modality Decomposition in Tokenizers / Detokenizers.** Due to the aligned data format and the unified observation and action spaces across embodiments, we are able to separately transform each semantically distinct component of the observational input and the action output, which we refer to as *modality*, and specify the compositional structure at the interface of the transformer trunk and the tokenizers / detokenizers. This design provides another layer of modularization to training and is core to the effectiveness of our method.

Concretely, for tokenization in the embodiment e, we encode the input observation  $o_t$  with multiple tokenizers  $\{T_{e,m_i}\}$  at the finer granularity of modalities denoted by  $o_t[m_i]$ . For instance, instead of aggregating all image inputs

before passing through the vision tokenizer, we use separate tokenizers for each camera view. All the encoded modalities are concatenated to compose the input tokens to the transformer trunk.

Similarly, for detokenization, we specify the subset of the transformer output tokens corresponding to each action modality, e.g. body pose, end effector pose, and gripper angles, and decode the selected tokens to yield each modality with separate detokenizers  $\{D_{e,m_i}\}$ . For convenience, we use the set of observation and action modalities as defined by the data collection formats in (3) and (4).

By explicitly decomposing the input and output modalities and encoding them separately, we are leveraging the innate structure of observations and actions and imposing such a structure on the token sequences processed by the transformer. Consequently, the knowledge of how to process different modalities learned during training can be shared across embodiments, fostering efficient transfer of the policy.

Although we employ a consistent data format and aligned input/output representations across embodiments, some modalities are not present or available for all embodiments. For example, the human operator is not equipped with a wrist camera, while the LocoMan robot has a wrist camera in some tasks to improve manipulation accuracy. In this case, we use masks defined during data collection to signify redundant dimensions in the observations as well as in the action labels. We refer the reader to Section IV-I for more implementation details.

In general, the highly modularized design of our learning framework offers great flexibility in handling all types of manipulation tasks across different embodiments, and effectively enhances the learning performance by capturing the common patterns in manipulation problems.

#### E. MXT Training Paradigm

We leverage the human data to pretrain the network for versatile manipulation policies. Specifically, for a given task, we first pretrain our network with the human dataset, and then finetune it with the LocoMan dataset (Algorithm 1). Only the transformer trunk weights are loaded from the pretrained checkpoint for finetuning. For certain tasks that are similar in nature but with different manipulation modes, we also collectively pretrain the model on the human datasets from these tasks, and then finetune on each task with the corresponding LocoMan dataset.

**Learning Objective.** We use the behavioral cloning objective for both pretraining and finetuning. In general, given a dataset  $\mathcal{D}_e$  on an embodiment e and aligned action modalities  $m_1, ..., m_k$ , the total loss to optimize when training on e is:

$$\mathcal{L}_{e}(\theta) = \sum_{i=1}^{k} \mathcal{L}_{e,m_{i}}(\theta),$$
(5)

where  $\mathcal{L}_{e,m_i}$  is the  $\ell_1$  loss of the action modality  $m_i$  with respect to the dataset of embodiment *e*. In practice, we optimize the following batched loss for each training batch

 $B_e = \{(\boldsymbol{o}_j, A_j)\}_{j=1}^n$  as a proxy of  $\mathcal{L}_{e,m_i}(\theta)$ :

$$\mathcal{L}_{e,m_{i}}(B_{e}) = \frac{1}{n} \sum_{j=1}^{n} \left[ \frac{1}{h} \sum_{l=1}^{h} \ell_{1} \left( \boldsymbol{a}_{j,l} \left[ m_{i} \right], \hat{\boldsymbol{a}}_{j,l} \left[ m_{i} \right] \right) \right], \quad (6)$$

where  $\mathbf{a}_{j,l}[m_i] = (A_j)_l[m_i]$  is the *l*-th step action of modality  $m_i$  in the action label sequence sample  $A_j = \{\mathbf{a}_{j,l}\}_{l=1}^{h}$ ;  $\hat{\mathbf{a}}_{j,l}[m_i] = [\pi_{\theta}(\mathbf{o}_j)]_l[m_i]$  is the predicted action of modality  $m_i$  at step *l*, and *h* is the chunk size or the action horizon.

Algorithm 1 Pretraining MXT on human data and finetuning on LocoMan data

**Require:** Human dataset  $\mathcal{D}_{human}$ , LocoMan dataset  $\mathcal{D}_{LocoMan}$ 

**Ensure:** Policy  $\pi$  for versatile LocoMan manipulation Initialize the MXT policy network  $\pi_{\theta}$  with parameters  $\theta$ . Set pretraining learning rate  $\eta_{\text{pretrain}}$ 

for step = 1, 2, ... do  $\triangleright$  Pretraining Stage Sample a batch *B* from  $\mathcal{D}_{human}$ 

Compute  $\mathcal{L}_{\text{human}}(B) = \sum_{i} \mathcal{L}_{\text{human},m_i}(B)$  with Eq.6 Optimize the policy weights  $\theta$  with backpropagation

Reinitialize the tokenizers and detokenizers of  $\pi$ . Preserve the trunk weights  $\theta_{trunk}$  learned from pretraining.

Set finetuning learning rate  $\eta_{\rm finetune}$ 

for step = 1, 2, ... do  $\triangleright$  Finetuning Stage Sample a batch *B* from  $\mathcal{D}_{\text{LocoMan}}$ 

Compute  $\mathcal{L}_{LocoMan}(B) = \sum_{i} \mathcal{L}_{LocoMan,m_i}(B)$  with Eq.6

Optimize the policy weights  $\theta$  with backpropagation **return**  $\pi$ 

# F. Experiment Setups

1) Tasks: We evaluate MXT on six household tasks of varying difficulty, across unimanual and bimanual manipulation modes of the LocoMan robot, with data collected by the Human2LocoMan system:

- Unimanual Toy Collection (TC-Uni). In this task, the robot must pick up a toy randomly positioned within a rectangular area and place it into a designated basket on the ground. Completing this task requires the robot to coordinate its whole-body motions to efficiently and accurately reach various locations on the ground and above the basket. As shown in Figure 5, we use 10 objects for robot finetuning and all objects for human pretraining and real-robot evaluation. The substeps of this task include: grasp the toy, and release the toy.
- *Bimanual Toy Collection (TC-Bi).* Similar to *Unimanual Toy Collection*, this task requires the robot to pick up a toy randomly placed within two rectangular areas on either side of a basket. We use 10 objects for robot finetuning, while all objects are included in human pretraining and real-robot evaluation. The substeps of this task include: grasp the toy, and release the toy.
- Unimanual Shoe Rack Organization (SO-Uni). This longer-horizon task involves organizing two shoes

placed on different levels of a shoe rack. The robot must coordinate whole-body motions to reach various rack levels and utilize both prehensile and non-prehensile manipulation skills. As shown in Figure 5, this task involves three pairs of shoes, with one pair being out-ofdistribution (OOD). The substeps of this task include: push the shoe on the higher rack, tap the shoe on the higher rack, transfer the gripper to the lower level, and tap the shoe on the lower rack.

- *Bimanual Shoe Rack Organization (SO-Bi).* One pair of shoes is randomly placed at the edge of the third level of the shoe rack. The robot must push one shoe inward and align it with the other. The substeps of this task include: push the shoe, and tap the shoe.
- Unimanual Scooping (Scoop-Uni). The robot performs unimanual manipulation using a litter shovel to scoop a 3D-printed cat litter from varying locations and poses within a litter box, and then dump it into a trash bin. This long-horizon task involves both tool use and deformable object manipulation. The task is decomposed into the following substeps: grasp the shovel, scoop the litter, tilt the shovel, dump the litter, and place the shovel back.
- *Bimanual Pouring (Pour-Bi).* The robot performs bimanual manipulation to pour a Ping Pong ball from one cup to another. This longer-horizon task requires the robot to accurately reach both cups, which are randomly placed within a rectangular area on a table, lift one cup, pour the ball into the other, and then place both cups back on the table. This task evaluates the coordination and precision of the robot's bimanual manipulation. The substeps of this task include: pick up both cups, pour the ball, and place both cups.

2) Human2LocoMan Embodiments: As shown in Table II, the unimanual and bimanual modes of Human2LocoMan represent distinct embodiments, each differing in morphology, observations, and action spaces. In practice, we install and utilize wrist cameras on the LocoMan robot for the three unimanual manipulation tasks.

TABLE II: Human2LocoMan embodiments (R=Right,L=Left).

Embodiments	Head Images	Wrist Image	Body Priop.	R-EEF Priop.	L-EEF Priop.	Body Pose	R-EEF Pose	L-EEF Pose	R-Grasp Action	L-Grasp Action
Human-Unimanual (R)	<b>↓</b> √	×	~	~	×	1	1	×	1	×
Human-Unimanual (L)	1	×	~	×	~	~	×	~	×	1
Human-Bimanual	1	×	~	~	~	~	~	~	~	~
LocoMan-Unimanual (R)	<ul><li>✓</li></ul>	~	~	~	√	~	~	×	~	×
LocoMan-Unimanual (L)	1	~	~	~	~	~	×	~	×	~
LocoMan-Bimanual	<ul> <li>✓</li> </ul>	×	1	~	1	×	~	~	~	~

3) Data collection: For each task, we collect various numbers of human and robot trajectories with the Human2LocoMan system. The details of the collected data are demonstrated in Table III. About 10% data of each task is used for validation.

4) Training details.: For Toy Collection and Shoe Rack Organzation, we pretrain a model that utilizes the human data of both the unimanual and bimanual versions of the task, then we finetune the model on each unimanual or bimanual task with the corresponding robot data. For each task, we



Fig. 5: Rollouts of the MXT policy and the objects used across manipulation tasks in our experiments. Green arrows indicate end-effector motions, red arrows denote torso movements, and pink arrows represent gripper actions. Both unimanual and bimanual toy collection tasks assess the robot's ability to grasp objects of varying shapes, colors, and positions. The unimanual variant emphasizes coordination between the torso and end-effector, while the bimanual variant highlights synchronized control of two loco-manipulators. Unimanual and bimanual shoe rack organization tasks evaluate non-prehensile manipulation skills such as pushing and tapping. The unimanual variant additionally requires torso articulation to reach shoes placed at different heights. Scooping is a complex task involving tool use, deformable object manipulators, and wide-range torso motion. Pouring is a long-horizon task that demands precise coordination of both loco-manipulators.

TABLE III: Records of data collection for different tasks.

Task	# human traj.	human time (min)	# robot traj.	robot time (min)
TC-Uni	300	25	150	15
TC-Bi	315	22	70	7
SO-Uni	240	34	90	23
SO-Bi	200	20	92	12
Scoop-Uni	340	96	66	22
Pour-Bi	210	35	64	22

choose a set of training hyperparameters (e.g. batch size, chunk size) that are kept the same for all methods. (See Section IV-K.) We also list the model hyperparameters we use for our method and the baselines in the Section IV-I and IV-J.

5) *Baselines:* We compare Human2LocoMan to the following SOTA imitation learning baselines:

- *Humanoid Imitation Transformer (HIT):* HIT [20] is an imitation learning framework designed for humanoid skill learning that also extends to any robot embodiment. It builds upon ACT [18] and employs a decoder-only architecture that simultaneously predict the future action sequence and future image features. It discourages the vision-based policy to ignore the visual input and overfit on proprioceptive states by introducing a L2 image feature loss to the original behavioral cloning policy. HIT itself is not capable of handling data from different domains and embodiments, and we position HIT as a reference implementation that efficiently learns from indomain robot demonstrations.
- *Heterogeneous Pretrained Transformer (HPT):* HPT [34] is a framework for learning from vast



Fig. 6: Substep success rate. The success rate for some substep is calcuated as the percentage of trials where the robot successfully completed the substep. For each task, we calculate this with 24 ID rollouts and 12 OOD rollouts. **MXT-Pretrained**: MXT pretrained on human dataset (including unimanual and bimanual if applicable), then finetuned on the LocoMan data. **MXT-Scratch**: MXT trained only on the LocoMan data. "L" denotes the larger training set (80 trajectories for SO-Uni, 60 trajectories for Pour and Scoop), while "S" denotes the smaller training set (40 trajectories for SO-Uni, 30 trajectories for Pour and Scoop).

amounts of data collected from humans, teleoperation, simulation, and real-life robots. HPT also has a modularized design and consists of the stems, the trunk, and the head, where the stems and heads are similar to our tokenizers and detokenizers. The trunk is designed to capture the complex mapping between the input and output in a unified latent space through large-scale pretraining. The implementation of HPT differs from our framework in several key aspects. Firstly, we leverage the unified observation and action frames to align data from different embodiments on the modality level, while HPT can only construct tokenizers for all image or proprioceptive data, and one detokenizer for all action dimensions. The ResNet image encoder in HPT is also frozen to achieve efficient learning with large models, while we opt to finetune the ResNet encoder along with the whole network end-to-end to better account for the visual gap between embodiments.

More implementation details of these baselines can be found in Section IV-J. For the HPT baseline, we train with several different settings: training with only LocoMan data, pretraining with our human data and finetuning on LocoMan data, and directly finetuning the released HPT checkpoints with LocoMan data. For the HIT baseline, we only train on LocoMan data, as it is unable to incorporate human data.



Fig. 7: Best validation loss of our method and HIT on all our tasks. **MXT-Pretrained**: MXT pretrained on human dataset (including unimanual and bimanual if applicable), then finetuned on the LocoMan data. **MXT-Scratch**: MXT trained only on the LocoMan data. The number suffix denotes the number of demonstrations in the LocoMan training set.

6) Evaluation Metrics: We present the evaluation results using three metrics: i) success rate (SR), ii) task score (TS), and iii) validation loss. To calculate the success rate and task score, we perform a fixed number of real world rollouts using the evaluated method for one task. The policy is rolled out for 24 times with in-distribution (ID) objects and 12 times with out-of-distribution (OOD) objects.

For each task, we define a set of critical substeps necessary to fully complete the task. When calculating the task score, successfully completing each intermediate substep earns one point, and reaching the final goal—i.e., completing the entire task—earns an additional point. The final task score is the sum of points across all rollouts for that task. The success rate of a method on a given task, under either the ID or OOD setting, is computed as the ratio of successful rollouts (i.e., rollouts where all substeps are completed) to the total number of rollouts performed.

In addition, we report the best validation loss as another metric for training performance. For all the included methods, we align how the loss is computed so that these losses can be meaningfully compared. Note that the validation loss is not a faithful indicator of the policy performance, but it does reflect how well the model is optimized, especially when there is a significant difference in the validation loss of different policies in the same setting. We mainly use this metric to analyze the training process of different architectures (MXT, HIT and HPT) and to provide a separate dimension to our evaluation.

## G. Supplementary Results and Analysis

(1) Does the Human2LocoMan system enable versatile quadrupedal manipulation capabilities? Data collection. As shown in Table III, Human2LocoMan teleoperation enables the collection of a substantial amount of robot data (over 50 trajectories) within 30 minutes across all tasks. Using the Human2LocoMan human data collection system, over 200 trajectories can be gathered within the same time frame. Even for the most challenging task, a human can collect over 300



Fig. 8: Best validation loss of our method and HPT on the unimanual Toy Collection task. **MXT-Pretrained**: MXT pretrained on human dataset (including unimanual and bimanual if applicable), then finetuned on the LocoMan data. **MXT-Scratch**: MXT trained only on the LocoMan data. **HPT-Pretrained**: HPT trunk pretrained on our human data, then finetuned on the LocoMan data. **HPT-Scratch**: HPT network trained only on the LocoMan data. **HPT-Base**: Finetune with our LocoMan data with HPT trunk initialized with released HPT-Base weights. **HPT-Small**: Finetune with our LocoMan data with HPT trunk initialized with released HPT-Small weights.

trajectories within one and a half hours. Notably, the robot's manipulation speed is comparable to, and in many tasks approaches, that of a human. These results highlight the data collection efficiency of our system. Task versatility. As depicted in Figure 5, Human2LocoMan's policy can perform tasks across a wide range of scenarios, including unimanual and bimanual manipulation, prehensile and non-prehensile manipulation, deformable object manipulation, and tool use, while also generalizing to OOD objects and conditions. Task performance. We summarize the success rates and task scores of our method and HIT across all tasks in Table I. Human2LocoMan's MXT achieves strong performance on all tasks using a relatively small dataset. The baseline method also attains decent performance on most tasks. These results highlight the high quality of our collected data and demonstrate the effectiveness of Human2LocoMan's data collection and training pipeline.

(3) How does human data collected by Human2LocoMan contribute to imitation learning performance? Long-horizon performance. For a more detailed analysis on long-horizon tasks that require multiple manipulation steps, we present in Figure 6 how the success rate decays with each substep in tasks including SO-Uni, Pour-Bi and Scoop-Uni. MXT-Pretrained is shown to maintain a decent success rate as the long-horizon task progresses, while MXT-Scratch and HIT tend to fail more after the first substep, especially in Pouring and Scooping tasks. We note that the second substep in these tasks commonly involves moving and localizing an object with precision, and pretraining with human data appears to help with completing such challenging substeps. This suggests that human data incorporated during pretraining can promote manipulation accuracy, which is key to completing a sequential long-horizon task.

## H. Related Work

Embodiments for Diverse Loco-Manipulation Skills: Learning manipulation skills on quadrupedal robots has shown promise and popularity in recent years, due to the versatility and mobility of the platforms. Many manipulator configurations and capabilities have been proposed for quadrupeds, including non-prehensile manipulation using the quadruped's legs or body (e.g., dribbling a soccer ball, pressing buttons, closing appliance doors, etc.) [35]–[42], using a back-mounted arm for tabletop tasks [8], [43], or using leg-mounted manipulators for spatially-constrained (e.g., reaching toys underneath furniture) or bi-manual manipulation tasks [14]. In this work, we take inspiration from the open-source LocoMan hardware platform [14], with two legmounted manipulators, which enable the training of policies across challenging tasks and multiple operating modes.

Learning Versatile Quadrupedal Manipulation: Reinforcement learning (RL) has been used for training individual non-prehensile manipulation skills [35], [36], [38]-[42], [44]-[49] and for training whole-body controllers to track end-effector poses for uni-manual grasping [8]-[10], [50]–[53]; here, policies are trained in simulation then transferred to the real robot platform, often with high cost in training complexity and training time. To mitigate some of these issues, imitation learning (IL) allows robots to directly learn from expert demonstrations [15], [54]-[56] and thus provides an alternative approach for efficiently acquiring more general manipulation skills [26], [57]-[60]. However, collecting robot data for quadrupedal platforms remains challenging, due to their high degrees of freedom and the need for stable whole-body controllers. Prior works have trained non-prehensile quadrupedal manipulation policies by learning from demonstrations collected in simulation [12], or grasping policies for a top-mounted arm using data collected from real-world demonstrations [10], [11], [13]. Our work introduces a scalable way of achieving more versatile manipulation skills on quadrupedal platforms encompassing both single-gripper and bi-manual manipulation tasks, using a small amount of robot data combined with human demonstrations collected via our novel teleoperation and data collection system.

**Data Collection for Imitation Learning:** Various methods have been utilized to collect data for imitation learning. Joysticks and spacemouses [16], [61], [62] are commonly used to directly teleoperate the robot for data collection. Cameras are employed to capture human motions and map them to the robot [17], [20], [63]–[65]. VR controllers provide a more intuitive way for the human to teleoperate the robot with visual or haptic feedback for dexterous manipulation tasks on robot arms, quadrupeds, and humanoid robots [13], [21], [22], [31], [66]–[68]. While most above works teleoperate the robot in task space, other works employ ex-skeleton or leader-follower systems to collect robot demonstrations by mapping the joint positions of the leader system to the robot [18], [19], [23], [31], [69]. To ease the burdens of teleoperating real robots and to scale up data collection,

recent works have achieved success by collecting human demonstrations in the wild with AR-assist [30] or handheld grippers [11], [70], although these are constrained to a specific robot or end-effector type. Other works enable more ergonomic data collection with body-worn cameras [27], [71] or VR glasses [31]. We introduce a unified framework to collect cross-embodiment data including both robot and human demonstrations, where the teleoperation system considers the whole-body motions of the embodiments to extend its workspace and actively sense the environment. The different manipulation modes of both the robot and human are regarded as different embodiments and the collected data can be used for model pre-training.

**Cross-Embodiment Learning:** Drawing from the success of foundation models in computer vision and natural language, there are many endeavors to replicate the success in robotics by training generalist policies on large-scale data from different embodiments [34], [72]–[76]. However, this remains an open challenge due to the heterogeneity of robot embodiments, and gaps in kinematics, vision, and proprioception.

Different neural architecture were proposed to handle the heterogeneity. CrossFormer [76] formulated policy learning as a sequence-to-sequence problem, so that any number of camera views or proprioceptive sensors can be handled as sequence of tokens, and add special readout tokens as part of the input sequence. In comparison, HPT [34] features a modularized structure and maps the variable observations to a fixed number of number tokens. In our work, we propose Modularized Cross-embodiment Transformer (MXT) that also employs a modularized design, but further enhances the modularity by identifying fine-granular alignment of data modalities between embodiments.

Notably, EgoMimic [31] proposed the idea that human be treated as another embodiment and demonstrated positive transfer by co-training on human and robot data. To achieve such positive transfer, EgoMimic minimizes human and robot kinematic gap by choosing a human-like robot embodiment, proprioception gap by normalizing and align action distributions, and appearance gap with visual masking. In comparison, Human2LocoMan is more flexible and scalable, transferring from human to quadruped without explicit domain alignment.

# I. Implementation and Training details of MXT

**Training Details.** We list the training optimizer and the transformer trunk hyperparameters in Table IV. These hyperparameters are kept the same for all our experiments.

**Cross Attention in Tokenizers and Detokenizers.** In the tokenizers of MXT, we use a simple cross attention mechanism to transform the input feature of indefinite length into a fixed number of tokens. For the attention layer in all tokenizers, the hidden dim is 128, the number of attention heads is 4, each with a head dimension of 32, and the dropout rate is 0.1. Other hyperparameters of each tokenizer are shown in Table V.

### TABLE IV: MXT trunk and training hyperparameters

Hyperparameters	Value
optimizer	AdamW
learning rate	5e-5 (finetuning/from scratch) 1e-4 (pretraining)
scheduler	constant
weight decay	1e-4
trunk encoder layers	4
trunk decoder layers	4
hidden dim	128
transformer feedforward dim	256
#attention heads	16

Similarly, we also use cross attention to decode the action modalities in detokenizers from a fixed number of output transformer tokens. For the attention layer, the number of attention heads is 4, each with a head dimension of 16, and the dropout rate is 0.1. Other hyperparameters of each detokenizer are shown in Table VI

TABLE V: MXT tokenizer hyperparameters

Modality	Input dimensions	#tokens	MLP widths
main image wrist image	(3, 480 1280) (3, 480, 640)	16 8	N/A
body pose EEF pose EEF to body pose gripper angles	(6,) (12,) (12,) (2,)	4 4 4	[128, 128]

TABLE V	VI: MXT	detokenizer	hyperparameters
---------	---------	-------------	-----------------

Modalities	Output dimensions	#tokens
body pose	(6,)	6
EEF pose	(12,)	6
gripper angle	(2,)	6

**Masks for aligning embodiment modalities.** We mentioned that masks are needed to exclude redundant dimensions or modalities that are not present in some embodiment, and here we give a more detailed description of our implemented masks.

a) Masks on images. We recognize that some image view are not available for all embodiments and tasks. In our current framework, we assume there are at most two camera views (or image modalities): the main camera and the wrist camera. However, this can be easily extended within our framework to cater to any number of camera views. When one of these camera views is not present, we directly mark this in the transformer mask of the trunk and fill in dummy tokens in the corresponding positions, so that the positions associated with this image modality will not be attended on.

b) Masks on proprioceptive states. In some cases, the proprioceptive states may have some or all dimensions that should not be considered for the task. For example, in singlearm tasks, the poses of the left end effector, or the last half of the end effector pose modality, will not be considered, and in bimanual tasks where the LocoMan body is upright, the body pose is fixed and therefore redundant in the observations. When part of a proprioception modality are redundant dimensions, we apply zero padding on these dimension and perform encoding through the tokenizer as usual. Different from how we treated masked image modalities, this has no effect on the transformer mask of the trunk. When an entire proprioception modality should be disregarded, however, we handle this modality in a similar to the image modalities and apply the transformer mask accordingly.

**Data Normalization.** For both human and LocoMan data, we apply data normalization on observations and action labels. For non-image data, we estimate the per-dimension mean the standard deviation from the dataset, and normalize the data with the usual approach:

$$\bar{x}_t = \frac{x_t - \text{mean}}{\text{std}}$$

For image data, the mean and standard deviation are set as the ImageNet statistics for the RGB channels: mean = [0.485, 0.456, 0.406], and std = [0.229, 0.224, 0.225]. The images are normalized in the same way with these parameters.

**Dropout in Pretraining.** We discover that increasing the dropout in transformer trunk can improve the finetuning performance for MXT. In practice, we find that setting the pretraining dropout rate to 0.5 for scooping and 0.4 for all other tasks yield reasonably good performance. When training with LocoMan data, including training from scratch and finetuning, the transformer trunk dropout rate is reverted to 0.1.

TADLE VII. HIT Hyperparameter	TABLE	VII:	HIT	hyperparameters
-------------------------------	-------	------	-----	-----------------

Hyperparameters	Value
optimizer	AdamW
learning rate	2e-5
scheduler	constant
weight decay	1e-4
encoder layers	4
decoder layers	4
hidden dim	128
#attention heads	8
feature loss weight	0.001
image backbone	ResNet18

## J. Implementation details of baselines

**HIT.** Our implementation of Humanoid Imitation Transformer [20] is based on the released codebase, with only minor modifications to accommodate our data format. The hyperparameters used for training are summarized in Table VII.

**HPT.** We follow the original implementation of HPT [34], with the main exception that we changed the data normalization method so as to align with the approach of other frameworks and to have a fair comparison of the validation loss. The hyperparameters we used when training HPT are summarized in Table VIII.

## TABLE VIII: HPT hyperparameters

Hyperparameters	Value
optimizer	AdamW
learning rate	5e-5 (finetuning/from scratch) 1e-4(pretraining)
scheduler	constant
weight decay	1e-4
trunk	
#transformer blocks	16
hidden dim	128
feedforward dim	256
#attention heads	8
action head	
#attention heads	8
head dim	64
dropout	0.1
output dim	20
image stem	
encoder	ResNet18
MLP widths	[128]
#tokens	16
state stem	
MLP widths	[128]
#tokens	16

## K. Global task-specific training parameters

We choose a set of training parameters for each specific task, and we keep these settings aligned across all methods as listed in Table IX.

TABLE IX: Global training parameters for each task

Mode	Batch Size	Training Steps	Chunk Size
Unimanual Bimanual	16 16	60000 60000	60 60
Unimanual Bimanual	24 24	80000 100000	180 120
Unimanual	24	100000	120
Bimanual	24	80000	180
	Mode Unimanual Bimanual Bimanual Unimanual Bimanual	ModeBatch SizeUnimanual16Bimanual24Unimanual24Unimanual24Bimanual24Bimanual24	ModeBatch SizeTraining StepsUnimanual1660000Bimanual1660000Unimanual2480000Bimanual24100000Unimanual24100000Bimanual2480000