# Improved Reinforcement Learning Coordinated Control of a Mobile Manipulator using Joint Clamping

Denis Hadjivelichkov[1], Kostas Vlachos[2], Dimitrios Kanoulas[1]

*Abstract*— Many robotic path planning problems are continuous, stochastic, and high-dimensional. The ability of a mobile manipulator to coordinate its base and manipulator in order to control its whole-body online is particularly challenging when self and environment collision avoidance is required. Reinforcement Learning techniques have the potential to solve such problems through their ability to generalise over environments. We study joint penalties and joint limits of a state-of-the-art mobile manipulator whole-body controller that uses LIDAR sensing for obstacle collision avoidance. We propose directions to improve the reinforcement learning method. Our agent achieves significantly higher success rates than the baseline in a goal-reaching environment and it can solve environments that require coordinated whole-body control which the baseline fails.

Fig. 1. The simulated mobile manipulator, composed by an omnidirectional mobile base and a 7DoF arm.

## I. INTRODUCTION

Mobile robots have a plethora of applications ranging from warehouse services, through oil rig inspections, to emergency interventions [1], [2], [3]. Modern robots require both high mobility and accurate manipulation to traverse collision-free paths while performing their tasks, which can be achieved via *mobile manipulators*. By applying Whole-Body Control (WBC), the base and manipulator movements of mobile robots coordinate to improve the efficiency of the system.

Classical WBC methods include the use of kinematic, velocity, and impedance controllers, model predictive controllers, and combinations thereof in advanced adaptive control strategies [4], [5], [6]. They have been shown to work well in many environments while also providing stability guarantees. Reinforcement Learning (RL) methods have shown great promise in their ability to compete with and potentially overcome classical methods in many robotic problems as they can work with complex inputs [7] and learn complex task solutions [8]. Once trained, RL agents can execute policies online, bringing down total mission times.

The recent state-of-the-art works on RL for mobile manipulator WBC have focused on goal reaching scenarios. While they show that their solutions are quicker than traditional methods they still underperform them in terms of success rate and impose big limitations on the robots and their environments, such as limited DoF [9] and simplistic tasks [10].

In this paper, we examine the trained WBC behaviour of a state-of-the-art baseline agent trained with a shaped reward. The robot agent is comprised of an omnidirectional base and a 7 DoF manipulator simulated in PyBullet (see Fig. 1). We determine potential causes of sub-optimality in the baseline
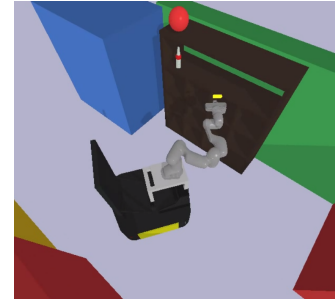
[1] University College London, denis.hadjivelichkov.19, d.kanoulas@ucl.ac.uk
[2] University of Ioannina, kostaswl@cse.uoi.gr

- frequent early episode termination due to reaching joint limits and consistently folded arm. We show that clamping the joint instead of a joint-limit penalty in the the reward improves the model's performance significantly and allows it to reach the goal much closer.

The summarised contributions of this paper are: (i) Identifying issues and potential improvements of a state-of-the-art method; (ii) Showing our method leads to higher success rates than the baseline and solves an environment requiring whole-body control, which the baseline fails; (iii) Evince our method's ability to generalise in an unseen environment.

## II. RELATED WORK

In this section, we present traditional and reinforcement learning approaches to whole-body control, justifying its use.

**Traditional Approaches:** Traditional approaches to implementing WBC include the use of kinematic and dynamic controllers [11], [12], [13]. Their advantage is that the current understanding of physical systems is refined and works well on fully actuated robots. Most methods focus on WBC for quadrupeds [14], [15], [16], humanoids [17], [18], [19], [20], [21], animaloids [22], [23], [24], or mobile manipulators [17], [25]. Model Predictive Control methods are popular with works such as Minniti et al. [26] showing success in WBC pose-tracking and interaction tasks. Recent works focus on non-linear strategies, such as Hierarchical Quadratic Programming [6], and non-linear model predictive control [5]. While some methods such as Operational Space Control can solve tasks with optimality and continuity in real-time, most traditional methods require large offline computation. Moreover, the methods for mobile manipulators are often based on simplified models of the robot which sometimes results in control solutions that are limiting the its agility.

**Reinforcement Learning Approaches:** Reinforcement learning approaches offer a framework that is transferable to

different tasks and robots, able to work online with scaling complexity, in a trade off with the limited prior information that it can use and long training that is often required for good performance. However, current methods still use application-specific architectures and rarely generalize to multi-task scenarios [27]. RL methods have successfully taught robots dexterous vision-based manipulation tasks [28], [29], [30], [31] and navigation tasks [32], [33].

Most research is also focused on legged robots [34], [35], [36], [37]. Wang et al. [10] integrate the state-of-the-art RL algorithms with visual perception for WBC and propose an efficient framework for decoupling of visual perception from control, which enables easier sim-to-real transfer. However, the used environment is simple, consisting of a table in front of a robot. Kindle et al. [9] use a Proximal Policy Optimization (PPO) based agent to train end-to-end whole-body control policies for obstacle avoidance and tested on a real mobile manipulator achieving state-of-the-art results. Their model makes use of Automatic Domain Randomization and Continuous Learning to guide the agent toward a solution in a custom reach-and-grasp environment. A hand-crafted reward function is defined with components for collision, joint limits, safety distance, optimal path following and time. These recent works show sub-optimal performance, worse than comparable traditional methods.

## III. BACKGROUND

We consider a standard RL framework, which includes an agent interacting with an environment via actions and observations. Environment rewards are fed into an RL learning algorithm, which optimises the agent's policy and thus creates a feedback loop. The problem focuses on goal-reaching environments in which a success is defined as the uninterrupted holding of the robot agent's end-effector within a given tolerance distance from the goal. The environments' state space consists of front and rear LIDAR scans, arm joint positions, arm joint and base velocities, and the goal location in the end-effector frame, while action space consists of joint and base accelerations. Both LIDAR observations and joint actions are limited to 2D planes. In this section, we discuss he state-of-the-art baseline [9] used for our experiments.

*1) Reward:* The baseline's reward function is handcrafted, encouraging the agent to learn to imitate a traditional path planning method and complete the task quicker, while discouraging it for moving close to objects. The reward has three termination cases: collision, timeout, and reaching joint limits. Finally, it also introduces an accumulation term that prevents the agent's exploitation of the reward.

*2) Agent:* The architecture of the agent is based on PPO with modified layers as depicted in Fig. 2. The two LIDAR scans are compressed via a separate scan block before being processed with the rest of the inputs in a network of fully connected layers. The agent produces a discretized policy for each action and its respective value.

## IV. METHOD

To understand the low success rate of the baseline in comparison with traditional methods, we analysed the behaviour
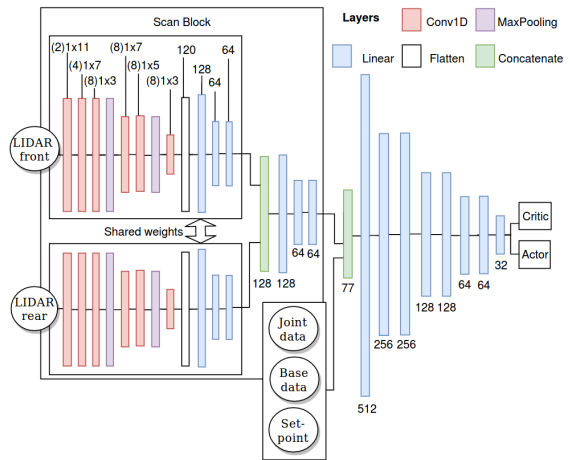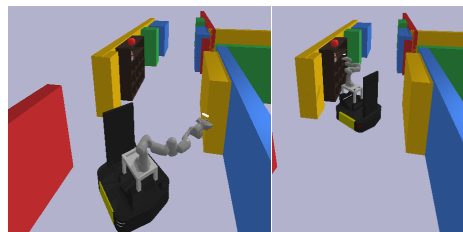


Fig. 2.  Agent's network architecture.



Fig. 3.  The robot controlled with the baseline agent at start (left) and end (right) of reaching task. Note the folded arm in the second image.

and performance of the agent after training. It was observed that the majority of episodes terminate due to the robot arm reaching its joint limits. We further noticed that this is a behaviour that can be limited explicitly instead of penalising and terminating the reinforcement learning agent.

While it was expected that the optimal solution would be for the robot to move toward the goal with a folded arm and unfold it while it is reaching the goal position, it was observed that the robot folds the arm in the beginning and does not change it throughout the whole run, as shown in Fig. 3 as well as the real robot experiments [38]. The cause of this could be partially explained by the custom environment itself, which does not explicitly require WBC in order to be solved. Additionally, the used reward function itself places more weight on optimal path following penalties than on timing penalties - following the optimal end-effector path is simpler when the manipulator is folded, because the end-effector is close to the base point of rotation. We attempt to address some of these drawbacks in this section.

**Improved Environments:** Two environments are used in our validations: a narrow corridor environment for comparison with state-of-the art and a new environment that cannot be solved without WBC.

The first environment, adapted from [9], consists of a narrow corridor of variable length, containing random avoidable obstacles and a randomly placed goal location (See Fig. 4). It requires the agent to plan its path and navigate around the obstacles toward the goal. We refer to this environment as *Corridor-env*.

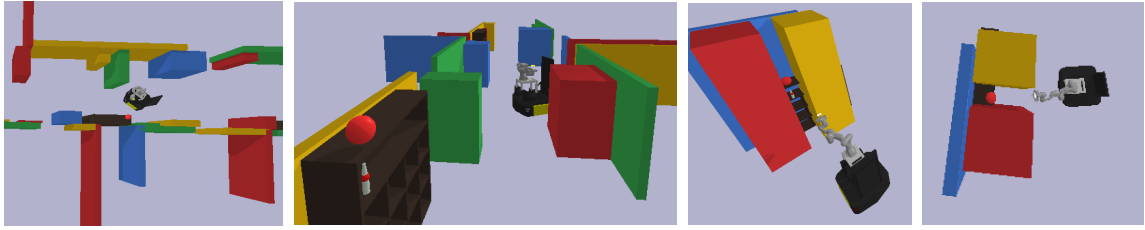However, its goal can be directly reached by folding

Fig. 4. The *Corridor-env* (left two) requires optimal path planning to be solved, but allows for the base and manipulator to move independently. The local *Gap-env* (right two) cannot be solved without coordinated control with the robot inserting its manipulator in the tunnel to approach the goal.

the arm and performing only mobile base collision-free navigation, thus does not require WBC.

We introduce a new environment (referred as *Gap-env*) which consists of narrow passages with the goal end-effector pose being reachable only with coordination between the base and manipulator (See Fig. 4). The width of the gap is very narrow and a small deviation of the arm or base during insertion would cause a collision, thus the task is difficult to solve without coordinated control. Two variants are used: In *Gap-env-train*, random uniform noise is added to initial joint angles, and orientation and position of the goal relative to the robot spawn location; In *Gap-env-test*, the tunnel gap width and length, as well as the goal placement relative to the tunnel is also randomly initialized to ensure that the testing scenarios are unseen by the agent.

In all the environments, a success is defined as the uninterrupted holding of the robot agent's end-effector within a given tolerance distance from the goal. Automatic Domain Randomization (ADR) is used to gradually adapt the complexity of the environment and guide the agent toward a solution. This is done by increasing or decreasing the acceptable tolerance distance to the goal depending on the agent's recent success rate. Via ADR, the tolerance distance to the goal is dynamically changed. The state space consists of 2D front and rear LIDAR scans, arm joint positions, arm joint and base velocities, and the goal location in end-effector frame. The action space is comprised of mobile base and arm joint accelerations. Given the complexity of the problem, its dimensionality is reduced to planar movements of the arm.

**Joint Clamping Method:** When training the baseline agent in the corridor environment, it was noticed that most of the episodes end due to joint limit termination. However, joint limits can easily be enforced by setting a manual limit (clamping) to the joint positions based on the robot's hardware limits with appropriate tolerance to protect the robot. Likewise, when training the baseline in the gap environment, it doesnt approach the goal, rather stays near the gap and oscillates. This is believed to be due to the baseline's safety margin penalty, which encourages the robot to keep its distance from all objects. We believe that a collision termination penalty is sufficient in teaching that behaviour. Thus, our modified reward function is as follows:

$$r_t = w_t \cdot \frac{\tau}{T_t} + w_{pd}\Delta d_{pd} + w_{pt}\frac{\Delta d_{pt}}{d_{pt,init}} + w_{ht}\frac{\tau}{T_h}$$
$$+ w_{hd}(1 - \min(1, d_g/d_h))\frac{\tau}{T_h} - I_h + D_c + D_h \quad (1)$$

where $w_t$ is the time penalty parameter, $\tau$ is the step time, and $T_t$ is the total time before episode timeout. This timeout reward encourages quicker task completion. An optimal path towards the goal is computed via Harmonic Potential Field (HPT). The goal distance reward penalises for deviation from the HPT $\Delta d_{pd}$ by $w_{pd}$ and rewards movement along the path $\Delta d_{pt}$ normalized by the total path $d_{pt,init}$ by $w_{pt}$. Furthermore, $w_{ht}$ is the reward for each time-step that the end-effector is within tolerance distance $d_h$ of the goal point and $w_{hd}$ is the reward for minimizing the distance to the goal position applied only when the distance to the goal $d_g$ is smaller than the tolerance distance $d_h$. Both of these rewards are normalized for the holding time threshold $T_h$ after which the task is done. $I_h$ is the accumulated holding reward which is subtracted if the end-effector leaves the tolerance sphere in order to prevent exploitation of the reward. Finally, $D_c$ is a collision penalty, and $D_h$ is the reward for sustained holding time $T_h$. The last two rewards end the current episode.

This reward function allows the agent to actuate the robot safely without hindering its learning and addresses the two issues encountered when running the baseline. The joint clamping is enforced programmatically, based on the robot's joint limits. The baseline agent is trained in with these modifications and compared with the standard baseline for several values of goal tolerance distance in Sec. V-A.

**Validation Setup** We use a mobile manipulator robot comprised of a omnidirectional base and a 7 DoF arm manipulator. All simulations are done using PyBullet 2.8 [39], while a high performance computing cluster is used for training. Agent parameters are shown in the appendix.

Our method and its compared baseline are trained in *Corridor-env* on 32 parallel workers for a total of 60M training steps. The final success rate, counted as number of successes over 100 episodes, is compared in Sec. V-A. We train the agents in *Gap-env-train* for 30M steps with 16 parallel workers. The agent is then evaluated in *Gap-env-test* and the results are reported in Sec. V-B. In both environments, the ADR is gradually adapting the tolerance distance in the range $[0.5; 0.05]$.

## V. RESULTS

We explore the following questions: (i) How does our method compare to the baseline? (ii) Is the new agent able to perform well on a task requiring Whole-Body Control? (iii) Is the agent generalizable to new environments?

Fig. 5. Our agent at start (left) and end (right) of reaching task.

| Tolerance Dist (m) | 0.5 | 0.2 | 0.1 | 0.07 |
|---|---|---|---|---|
| Baseline | 72% | 65% | 40% | 0% |
| Joint-Clamping Method (ours) | **73%** | **73%** | **63%** | **24%** |

TABLE I

SUCCESS RATES IN *Corridor-env* AGAINST TOLERANCE DISTANCES

| Environment | Gap-env-train | Gap-env-test |
|---|---|---|
| Baseline | fail | fail |
| Joint Clamping Method (ours) | 81% | 76% |

TABLE II

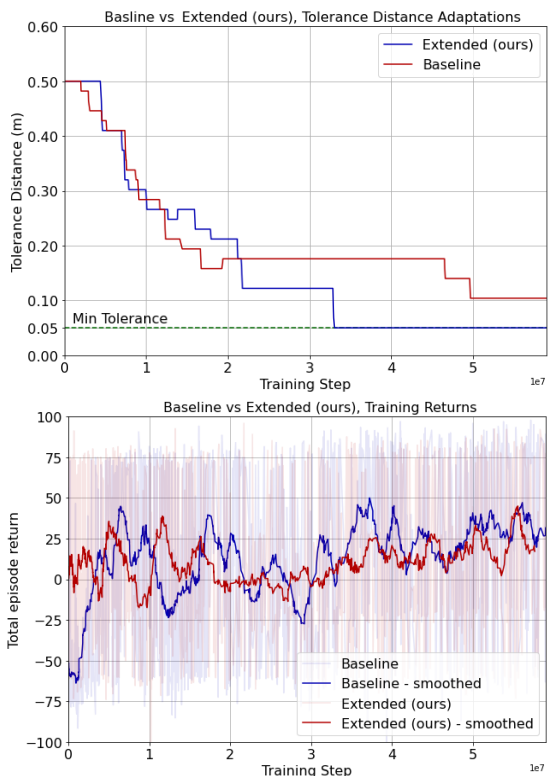SUCCESS RATES IN *Gap-env* WITH $0.05m$ TOLERANCE DISTANCE.



Fig. 6. **(left)** Adaptation of tolerance distance per step for both reward settings. **(right)** Total returns per episode plotted against the episode termination step. Running mean smoothing of 0.95 is used. Note that due to the use of ADR, the training rewards are fairly similar

### A. Validation on Corridor Environment

The original baseline agent with a joint limit penalty and our modified agent with clamped joint limits (Eq. 1) are ran with 32 workers for 60M steps. This process took 48 hours to finish. The trained models were tested in *Corridor-env* with fixed goal tolerance distances in the range of 0.5 to 0.07.

Running the baseline, we managed to achieve 72% for the highest tolerance distance (0.5m), while the success rate was significantly decreasing with the tolerance distance

dropping. The minimum successful tolerance distance was 0.1m. The resulting success rates of our agent, shown in the lower row of Table I, are noticeably higher than the baseline performance with standard reward, especially when the tolerance distance is decreasing.

This difference in performance can be further explained by the difference of ADR tolerances shown in Fig. 6. For the baseline, the ADR tolerance distance reached at 60M steps is 0.1, not reaching the lowest distance of 0.05. In comparison, our agent successfully adapted to the lowest tolerance distance 20M steps before the training ended. This indicates that with the modified reward, the agent learns quicker and better. The training returns of the original and modified baselines are shown in Fig. 6. For the modified agent in the environment with tolerance distance fixed to $0.07m$, it is found that 52% of the unsuccessful episodes terminate due to collision, while 48% terminate due to timeout. The distance from the end-effector to the goal at the end of unsuccessful episodes is on average $0.11m$. While the joint limit modification shows an increase in success rate, the folded manipulator behaviour is still observed (see Fig. 5).

### B. Whole-Body Control Task

Training the agent locally in *Gap-env*, which includes narrow tunnels where only the arm can fit, forces the simultaneous coordination between the arm and the mobile base as a WBC. With a 0.05m tolerance distance to the goal, the success rate of the training is 81%, shown in Table II. As can be observed from the success rates in the table, our agent successfully generalises to unseen variants of the training environment, with a drop of only 5% in success rate. In a typical episode, the robot is observed moving towards the goal, while adjusting its manipulator for tunnel-entry, as expected from a WBC solution. In failed episodes, it is seen that the robot often reaches the goal within less than $0.05m$, however it backs off and re-approaches several times until the episode terminates due to timeout. Note that the original baseline method was not able to solve such environments based on its handcrafted reward function.

## VI. CONCLUSIONS AND FUTURE WORK

This work presents an RL method for Whole-body Control of a mobile manipulator that improves on the state of the art. Our shaped reward function combined with joint limit clamping shows a significant improvement of 24% over the baseline for small tolerance distances. Moreover, the proposed agent is able to solve whole-body control tasks which the baseline fails. We show that training the agent with our reward in one environment, transfers its learned skills well to a similar, but different, testing environment.

The current method is limited to using 2D LIDAR data and planar manipulator actions. Future work will focus on expanding these limitations via more informative observations, such as 3D LIDAR or RGB-D images. More importantly, the "folding arm" behaviour should be further examined.

## REFERENCES

[1] S. Rajana, "Robotics for the Supply Chain," in *2017 Third International Conference on Science Technology Engineering & Management (ICONSTEM)*, 04 2018.

[2] M. Bengel, K. Pfeiffer, B. Graf, A. Bubeck, and A. Verl, "Mobile Robots for Offshore Inspection and Manipulation," in *2009 IEEE/RSJ International Conference on Intelligent Robots and Systems*, 12 2009, pp. 3317–3322.

[3] D. Rehak, A. Dudáček, and P. Poledňák, "A multipurpose robotic vehicle for the rescue of persons and interventions in emergency situations," *Komunikacie*, vol. 15, pp. 103–109, 01 2013.

[4] F. L. Moro and L. Sentis, *Whole-Body Control of Humanoid Robots*. Springer Netherlands, 2018, pp. 1–23. [Online]. Available: https://doi.org/10.1007/978-94-007-7194-9_51-2

[5] M. Logothetis, G. C. Karras, S. Heshmati-Alamdari, P. Vlantis, and K. J. Kyriakopoulos, "A Model Predictive Control Approach for Vision-Based Object Grasping via Mobile Manipulator," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2018, pp. 1–6.

[6] S. Kim, K. Jang, S. Park, Y. Lee, S. Y. Lee, and J. Park, "Whole-body Control of Non-holonomic Mobile Manipulator Based on Hierarchical Quadratic Programming and Continuous Task Transition," in *2019 IEEE 4th International Conference on Advanced Robotics and Mechatronics (ICARM)*, 2019, pp. 414–419.

[7] M. Jaderberg, V. Mnih, W. Czarnecki, T. Schaul, J. Leibo, D. Silver, and K. Kavukcuoglu, "Reinforcement Learning with Unsupervised Auxiliary Tasks," 2016.

[8] S. Levine, "Reinforcement Learning and Control as Probabilistic Inference: Tutorial and Review," [Online] Available: arXiv:1805.00909, 2018.

[9] J. Kindle, F. Furrer, T. Novkovic, J. J. Chung, R. Siegwart, and J. Nieto, "Whole-Body Control of a Mobile Manipulator using End-to-End Reinforcement Learning," [Online] Available: arXiv:2003.02637, 2020.

[10] C. Wang, Q. Zhang, Q. Tian, S. Li, X. Wang, D. Lane, Y. Petillot, and S. Wang, "Learning mobile manipulation through deep reinforcement learning," *Sensors*, vol. 20, no. 3, 2020. [Online]. Available: https://www.mdpi.com/1424-8220/20/3/939

[11] O. Khatib, "A unified approach for motion and force control of robot manipulators: The operational space formulation," *IEEE Journal on Robotics and Automation*, vol. 3, no. 1, pp. 43–53, 1987.

[12] M. Bauza and A. Rodriguez, "A probabilistic data-driven model for planar pushing," [Online] Available: arXiv:1707.06887, 2017.

[13] Y. Wu, P. Balatti, M. Lorenzini, F. Zhao, W. Kim, and A. Ajoudani, "A Teleoperation Interface for Loco-Manipulation Control of Mobile Collaborative Robotic Assistant," *IEEE Robotics and Automation Letters*, vol. 4, no. 4, pp. 3593–3600, 2019.

[14] G. Raiola, E. Mingo Hoffman, M. Focchi, N. Tsagarakis, and C. Semini, "A Simple Yet Effective Whole-Body Locomotion Framework for Quadruped Robots," *Frontiers in Robotics and AI*, vol. 7, p. 159, 2020.

[15] C. Dario Bellicoso, F. Jenelten, P. Fankhauser, C. Gehring, J. Hwangbo, and M. Hutter, "Dynamic locomotion and whole-body control for quadrupedal robots," in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2017, pp. 3359–3365.

[16] V. Morlando, A. Teimoorzadeh, and F. Ruggiero, "Whole-body control with disturbance rejection through a momentum-based observer for quadruped robots," *Mechanism and Machine Theory*, vol. 164, p. 104412, 2021.

[17] A. Dietrich, T. Wimböck, A. Albu-Schäeffer, and G. Hirzinger, "Reactive Whole-Body Control: Dynamic Mobile Manipulation Using a Large Number of Actuated Degrees of Freedom," *IEEE Robotics & Automation Magazine*, vol. 19, pp. 20–33, 2012.

[18] N. Mansard, O. Stasse, P. Evrard, and A. Kheddar, "A versatile Generalized Inverted Kinematics implementation for collaborative working humanoid robots: The Stack Of Tasks," in *2009 International Conference on Advanced Robotics*, 2009, pp. 1–6.

[19] E. M. Hoffman, B. Clément, C. Zhou, N. G. Tsagarakis, J.-B. Mouret, and S. Ivaldi, "Whole-Body Compliant Control of iCub: first results with OpenSoT," in *IEEE/RAS ICRA Workshop on Dynamic Legged Locomotion in Realistic Terrains*, Brisbane, Australia, 2018. [Online]. Available: https://hal.archives-ouvertes.fr/hal-01790597

[20] A. Rocchi, E. M. Hoffman, D. G. Caldwell, and N. G. Tsagarakis, "OpenSoT: A whole-body control library for the compliant humanoid robot COMAN," in *2015 IEEE International Conference on Robotics and Automation (ICRA)*, 2015, pp. 6248–6253.

[21] A. Laurenzi, D. Kanoulas, E. Mingo Hoffman, L. Muratore, and N. G. Tsagarakis, "Whole-Body Stabilization for Visual-Based Box Lifting with the COMAN+ Robot," in *Third IEEE International Conference on Robotic Computing (IRC)*, 2019, pp. 445–446.

[22] E.-J. Rolley-Parnell, D. Kanoulas, A. Laurenzi, B. Delhaisse, L. Rozo, D. G. Caldwell, and N. G. Tsagarakis, "Bi-Manual Articulated Robot Teleoperation using an External RGB-D Range Sensor," in *15th International Conference on Control, Automation, Robotics and Vision (ICARCV)*, 2018, pp. 298–304.

[23] V. S. Raghavan, D. Kanoulas, A. Laurenzi, D. G. Caldwell, and N. G. Tsagarakis, "Variable Configuration Planner for Legged-Rolling Obstacle Negotiation Locomotion: Application on the CENTAURO Robot," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2019, pp. 4738–4745.

[24] V. S. Raghavan, D. Kanoulas, D. G. Caldwell, and N. G. Tsagarakis, "Agile Legged-Wheeled Reconfigurable Navigation Planner Applied on the CENTAURO Robot," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2020, pp. 1424–1430.

[25] J. Liu, P. Balatti, K. Ellis, D. Hadjivelichkov, D. Stoyanov, A. Ajoudani, and D. Kanoulas, "Garbage Collection and Sorting with a Mobile Manipulatorusing Deep Learning and Whole-Body Control," in *20th IEEE-RAS International Conference on Humanoid Robots (Humanoids)*, 2020.

[26] M. V. Minniti, F. Farshidian, R. Grandia, and M. Hutter, "Whole-Body MPC for a Dynamically Stable Mobile Manipulator," *IEEE Robotics and Automation Letters*, vol. 4, no. 4, pp. 3687–3694, 2019.

[27] C. Sammut, "When do robots have to think?" *Advances in Cognitive Systems*, vol. 1, pp. 73–81, 2012.

[28] O. M. Andrychowicz, B. Baker, M. Chociej, R. Józefowicz, B. McGrew, J. Pachocki, A. Petron, M. Plappert, G. Powell, A. Ray, J. Schneider, S. Sidor, J. Tobin, P. Welinder, L. Weng, and W. Zaremba, "Learning dexterous in-hand manipulation," *The International Journal of Robotics Research*, vol. 39, no. 1, pp. 3–20, 2020. [Online]. Available: https://doi.org/10.1177/0278364919887447

[29] OpenAI, I. Akkaya, M. Andrychowicz, M. Chociej, M. Litwin, B. McGrew, A. Petron, A. Paino, M. Plappert, G. Powell, R. Ribas, J. Schneider, N. Tezak, J. Tworek, P. Welinder, L. Weng, Q. Yuan, W. Zaremba, and L. Zhang, "Solving rubik's cube with a robot hand," [Online] Available: arXiv:1910.07113, 2019.

[30] D. Kalashnikov, A. Irpan, P. Pastor, J. Ibarz, A. Herzog, E. Jang, D. Quillen, E. Holly, M. Kalakrishnan, V. Vanhoucke, and S. Levine, "QT-Opt: Scalable Deep Reinforcement Learning for Vision-Based Robotic Manipulation," [Online] Available: arXiv:1806.10293, 2018.

[31] R. Julian, B. Swanson, G. S. Sukhatme, S. Levine, C. Finn, and K. Hausman, "Efficient Adaptation for End-to-End Vision-Based Robotic Manipulation," [Online] Available: arXiv:2004.10190, 2020.

[32] T. Hester, M. Quinlan, and P. Stone, "RTMBA: A real-time model-based reinforcement learning architecture for robot control," in *2012 IEEE International Conference on Robotics and Automation*, 2012, pp. 85–90.

[33] A. Francis, A. Faust, H.-T. L. Chiang, J. Hsu, J. C. Kew, M. Fiser, and T.-W. E. Lee, "Long-range indoor navigation with PRM-RL," [Online] Available: arXiv:1902.09458, 2019.

[34] Z. Li, T. Zhao, F. Chen, Y. Hu, C. Su, and T. Fukuda, "Reinforcement Learning of Manipulation and Grasping Using Dynamical Movement Primitives for a Humanoidlike Mobile Manipulator," *IEEE/ASME Transactions on Mechatronics*, vol. 23, no. 1, pp. 121–131, 2018.

[35] J. Hwangbo, J. Lee, A. Dosovitskiy, D. Bellicoso, V. Tsounis, V. Koltun, and M. Hutter, "Learning agile and dynamic motor skills for legged robots," *Science Robotics*, vol. 4, no. 26, 2019. [Online]. Available: https://robotics.sciencemag.org/content/4/26/eaau5872

[36] C. Yang, K. Yuan, W. Merkt, T. Komura, S. Vijayakumar, and Z. Li, "Learning Whole-Body Motor Skills for Humanoids," in *IEEE-RAS International Conference on Humanoid Robots (Humanoids 2018)*, 2018, pp. 270–276.

[37] R. Lober, V. Padois, and O. Sigaud, "Efficient Reinforcement Learning for Humanoid Whole-Body Control," in *IEEE-RAS International Conference on Humanoid Robots*, Cancun, Mexico, 2016. [Online]. Available: https://hal.archives-ouvertes.fr/hal-01377831

[38] aslteam, "Whole-Body Control of a Mobile Manipulator using End-to-End Reinforcement Learning," https://www.youtube.com/watch?v=3qobNCMUMV4, [Online, accessed 22 August 2020].

[39] E. Coumans and Y. Bai, "PyBullet, a Python module for physics simulation for games, robotics and machine learning," http://pybullet.org, 2016–2020.

APPENDIX: HYPERPARAMETERS

| Parameter | Value | Parameter | Value | Parameter | Value |
|-----------|-------|-----------|-------|-----------|-------|
| $w_t$ | -15 | $d_h$ | 0.3 m | clip range | 0.2 |
| $w_{hd}$ | 40 | $D_h$ | 10 | clip range vf | -1 |
| $w_{ht}$ | 20 | $D_{jl}{}^*$ | -20 | noptepochs | 30 |
| $w_{pd}$ | -10 | $D_c$ | -60 | gamma | 0.999 |
| $w_{pt}$ | 50 | $\tau$ | 0.04 s | n steps | 2048 |
| $w_{sm}{}^*$ | -1 | $T_h$ | 1 s | nminibatches | 8 |

TABLE III

REWARD PARAMETERS USED FOR THE REWARD AND PPO AGENT.

PARAMETERS MARKED WITH $*$ ARE ONLY USED FOR THE BASELINE