# Deep Reinforcement Learning for Bipedal Locomotion: A Brief Survey

Lingfan Bao[1], Joseph Humphreys[1], Tianhu Peng[1] and Chengxu Zhou[2]

*Abstract*—Bipedal robots are garnering increasing global attention due to their potential applications and advancements in artificial intelligence, particularly in Deep Reinforcement Learning (DRL). While DRL has driven significant progress in bipedal locomotion, developing a comprehensive and unified framework capable of adeptly performing a wide range of tasks remains a challenge. This survey systematically categorizes, compares, and summarizes existing DRL frameworks for bipedal locomotion, organizing them into end-to-end and hierarchical control schemes. End-to-end frameworks are assessed based on their learning approaches, whereas hierarchical frameworks are dissected into layers that utilize either learning-based methods or traditional model-based approaches. This survey provides a detailed analysis of the composition, capabilities, strengths, and limitations of each framework type. Furthermore, we identify critical research gaps and propose future directions aimed at achieving a more integrated and efficient framework for bipedal locomotion, with potential broad applications in everyday life.

*Index Terms*—Deep Reinforcement Learning, Humanoid Robots, Bipedal Locomotion, Legged Robots

Fig. 1: **Common bipedal and humanoid robots used as platforms for testing DRL frameworks**. (a) NAO, a toy-like 3D humanoid robot, actuated by servo motors [5]. (b) Rabbit, a 2D bipedal robot, actuated by torque control [6]. (c) Cassie, a 3D bipedal robot, also actuated by torque control [7]. (d) ATLAS, a 3D humanoid robot, driven by hydraulics [8]. (e) Digit, a full human-sized 3D humanoid robot, an upgrade based on Cassie and actuated by torque control [9].

## I. INTRODUCTION

Humans navigate complex and varied environments, performing diverse locomotion tasks with only two legs. To facilitate dynamic bipedal locomotion, model-based methods were introduced in the 1980s and have since evolved significantly [1], [2], [3]. These methods, characterized by rapid convergence, furnish a predictive framework for understanding environmental structures. However, they struggle to adapt in dynamically challenging environments that are difficult to model precisely. More recently, advancements in machine learning have provided novel approaches. Reinforcement learning (RL)-based methods, in particular, are adept at navigating the full dynamics of robot-environment interactions [4]. Additionally, hybrid approaches that combine model-based and learning-based methods have been developed to leverage the advantages of both. Yet, the question remains: *Is there a unified framework capable of enabling bipedal robots to effectively manage a diverse range of locomotion tasks?*

To address this, we explore recent advancements in deep reinforcement learning (DRL)-based frameworks, categorizing control schemes into two primary types: (i) end-to-end and (ii) hierarchical. End-to-end frameworks map robot states directly to control outputs at the joint level, while hierarchical frameworks adopt a structured approach, decomposing
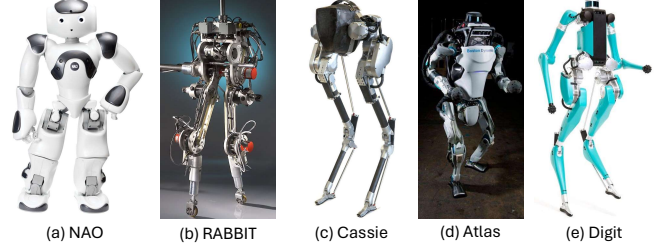
decision-making into multiple layers. Here, a High-Level (HL) planner addresses navigation and path planning, while a Low-Level (LL) controller focuses on fundamental locomotion skills. The highest decision-making tier, the task level, receives direct input from task or user commands.

The evolution of RL in bipedal robotics has spurred a dynamic growth in innovative applications. Although the application of RL to simple 2D bipedal robots began in 2004 [10], [11], it took several years before deep reinforcement learning (DRL) algorithms emerged. These DRL-based methods have since shown promising results in physical simulators [12], [13], [14]. Agility Robotics introduced the first sim-to-real end-to-end learning framework in 2019, which was applied on the 3D torque-controlled bipedal robot Cassie, as shown in Fig. 1(c) [7]. Besides model-based reference learning, the policy can also incorporate motion capture data [15], [16], [17], [18], or start from scratch [19] to explore solutions freely. Recent studies demonstrate that end-to-end frameworks robustly handle complex and diverse tasks [20], [21], [22].

Similarly, hierarchical structures have garnered significant interest. Within this subset, the hybrid approach combines RL-based and model-based methods to enhance both planning and control strategies. A notable framework employs a learned High-Level (HL) planner coupled with a Low-Level (LL) model-based controller, often referred to as the Cascade-structure or Deep Planning Hybrid Scheme [23], [24], [25]. Another innovative construction integrates a learned feedback controller with an HL planner, categorized under the DRL Feedback Control Hybrid Scheme [26], [27]. Additionally, a

[1]School of Mechanical Engineering, University of Leeds, UK. {`mnlb, el20jeh, mntp`}`@leeds.ac.uk`

[2]Department of Computer Science, University College London, UK. `chengxu.zhou@ucl.ac.uk`
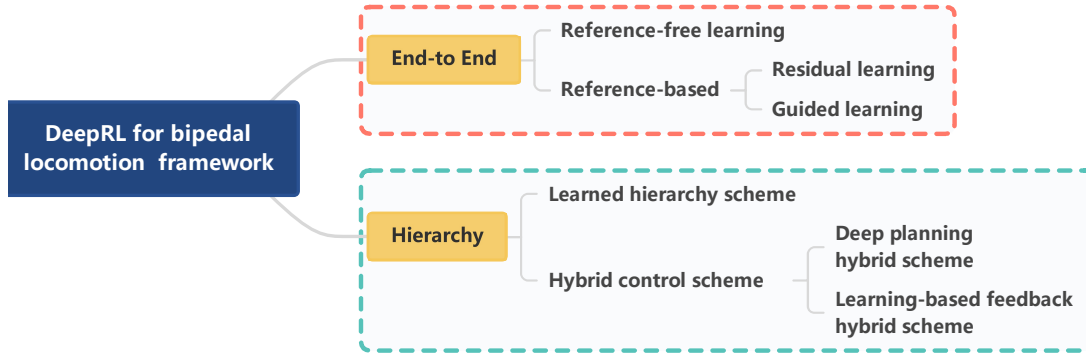
Fig. 2: **Classification of DRL-based control schemes.**

learned hierarchical control scheme [28] decomposes locomotion into various tasks, focusing each layer on specific functions such as navigation and fundamental locomotion skills [12], [13], [29].

While several review papers discuss RL for general robotics [4] and model-based methods for bipeds [1], [2], [3], none specifically focus on DRL-based frameworks for bipeds. This survey aims to address this gap by summarizing current research progress, highlighting the structure and capabilities of bipedal locomotion frameworks, and exploring future directions. We also catalogue DRL-based frameworks, as depicted in Fig. 2. The primary contributions of this survey are:

- A comprehensive summary and cataloguing of DRL-based frameworks for bipedal locomotion.
- A detailed comparison of each control scheme, highlighting their strengths, limitations, and characteristics.
- The identification of current challenges and the provision of insightful future research directions.

The paper is organized as follows: Section II focuses on end-to-end frameworks, categorized by learning approaches. Section III details hierarchical frameworks, classified into three main types. Section IV addresses existing gaps, ongoing challenges, and potential future research directions. Finally, Section V concludes the paper.

## II. END-TO-END FRAMEWORK

The end-to-end DRL framework represents a holistic approach where a single neural network (NN) policy, denoted $\pi(\cdot) : \mathcal{X} \rightarrow \mathcal{U}$, directly translates sensory inputs—such as images, lidar data, or proprioceptive feedback [30]—along with user commands [19] or pre-defined references [31], into joint-level control actions. These actions encompass motor torques [32], positions, and velocities [15]. This framework obviates the need for manually decomposing the problem into sub-tasks, streamlining the control process.

End-to-end strategies primarily simplify the design of low-level tracking to basic elements, such as a Proportional-Derivative (PD) controller. These methods can be broadly categorized based on their reliance on prior knowledge into two types: reference-based and reference-free. The locomotion skills developed through these diverse learning approaches exhibit considerable variation in performance and adaptability.

In the following sections, we will delve into various representation frameworks, exploring their characteristics, limitations, and strengths in comprehensive detail. To facilitate an understanding of these distinctions, Table I provides a succinct overview of the frameworks discussed.
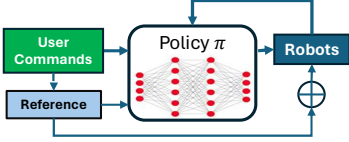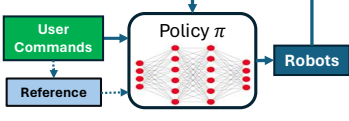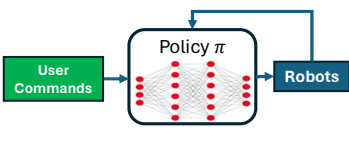
### A. Reference-based learning

Reference-based learning utilizes prior knowledge, allowing the policy to develop locomotion skills by adhering to predefined references, which may be derived from trajectory optimization (TO) techniques or captured through motion capture systems. This method facilitates the acquisition of locomotion skills compared to alternative approaches, though it typically results in locomotion patterns that closely resemble the predefined references or motion clips, thus limiting the variety of gait patterns. Generally, this approach can be divided into two primary methods: (i) residual learning and (ii) guided learning.

*1) Residual learning:* This method involves a framework that is aware of the current reference joint positions and applies offsets determined by the policy to modify motor commands at the current timestep. By utilizing predefined motion trajectories, the residual term acts as feedback control, compensating for errors and enabling the biped to achieve dynamic locomotion skills.

Introduced in 2018, a residual learning framework for the bipedal robot Cassie marked a significant advancement [33]. This framework allowed the robot to walk forward by incorporating a policy trained via Proximal Policy Optimization (PPO) algorithms, as detailed in Appendix A. The policy receives the robot's states and reference inputs, outputting a residual term that augments the reference at the current timestep. These modified references are then processed by a Proportional-Derivative (PD) controller to set the desired joint positions. While this framework enhanced the robot's ability to perform tasks beyond standing [39], its physical deployment on a bipedal robot has not yet occurred, potentially rendering it impractical for managing walking at varying speeds and limiting movement to a single direction.

To transition this framework to a real robot, a sim-to-real strategy based on the previous model was demonstrated, where the policy, trained through a residual learning approach, was subsequently applied on a physical bipedal robot [34]. This process and its key points are further explored in Appendix

TABLE I: Summary and comparison of reference-based and reference-free Learning approaches for end-to-end framework. The dashed line in the implementation flow chat refers to optional.

| Methods | Works | Capabilities | Characteristic | Advantages and Disadvantages | Implementation Flow Chart |
|---|---|---|---|---|---|
| Residual learning | [33] [34] [35] | Forward walk undirectional walk Omni-walk | Add residential term to the known motor positions at the current time step. | **A:** Fast convergence speed **D:** Require high-quality predefined reference, limit to specific motions, and lack robustness to complicated terrains. |  |
| Guided learning | [36] [31] [37] [22] | Forward walk Versatile walk Versatile jump Versatile motions | Mimic the predefined reference and directly specifies joint-level command. | **A:** Accelerate the learning process and robust to the terrains. **D:** Limited to the predefined motions and lack adaptability to unforeseen changes in environment. |  |
| Reference-free learning | [19] [38] [20] | Periodic motions Stepping stones Vision-based | Learn locomotion skills from scratch without any prior knowledge. | **A:** High potential for gait exploration, high robust to the complicated terrain. **D:** Requires intensive reward shaping for gait patterns and relatively expensive computational training resources. |  |

**Forward Walk** involves the bipeds walking straight ahead. **Unidirectional Walk** enables the bipeds to move both forward and backward within a range of desired velocities. Omni-Walk grants the bipeds the ability to walk in any direction. **Versatile Walk** allows the bipeds to walk forward, backward, turn, and move sideways, providing extensive movement capabilities. **Periodic Motions** entails the execution of various repeated gait patterns, such as walking, hopping, or galloping. **Versatile Jump** refers to jump towards different desired targets. **Versatile Motions** cover performing a broad array of motions, both periodic and aperiodic such as jumping.

B. Compared to model-based methods, this training policy achieves faster running speeds on the same platform, underlining the considerable potential of DRL-based frameworks. However, the robot's movements remain constrained to merely walking forward or backward. A novel approach in residual learning was introduced to enable unidirectional walking, where the policy outputs a residual term added to the current positional states, facilitating gradual omnidirectional walking [35].

*2) **Guided learning**:* Guided learning trains policies to directly output the desired joint-level commands, eschewing the addition of a residual term. The reward structure in this approach is focused on closely imitating predefined references.

A sim-to-real framework that employs periodic references to initiate the training phase was proposed in [36]. In this framework, the action space directly maps to the joint angles, and desired joint positions are managed by joint PD controllers. The framework also incorporates a Long Short-Term Memory (LSTM) network, as detailed in Appendix A, which is synchronised with periodic time inputs. However, this model is limited to a single locomotion goal: forward walking. A more diverse and robust walking DRL framework that includes a Hybrid Zero Dynamics (HZD) gait library was demonstrated [31], achieving a significant advancement by enabling a single end-to-end policy to facilitate walking, turning, and squatting.

Despite these advancements, the parameterization of reference motions introduces constraints that limit the flexibility of the learning process and the policy's response to disturbances. To broaden the capabilities of guided learning policies, a framework capable of handling multiple targets, including jumping, was developed [37]. This approach introduced a novel policy structure that integrates long-term input/output (I/O) encoding, complemented by a multi-stage training

methodology that enables the execution of complex jumping maneuvers. An adversarial motion priors approach, employing a style reward mechanism, was also introduced to facilitate the acquisition of user-specified gait behaviors [18]. This method improves the training of high-dimensional simulated agents by replacing complex hand-designed reward functions with more intuitive controls.

While previous works primarily focused on specific locomotion skills, a unified framework that accommodates both periodic and non-periodic motions was further developed [22] based on the foundational work in [37]. This framework enhances the learning process by incorporating a wide range of locomotion skills and introducing a dual I/O history approach, marking a significant breakthrough in creating a robust, versatile, and dynamic end-to-end framework. However, experimental results indicate that the precision of locomotion features, such as velocity tracking, remains suboptimal.

Guided learning methods expedite the learning process by leveraging expert knowledge and demonstrating the capacity to achieve versatile and robust locomotion skills. Through the comprehensive evaluation [22], it is demonstrated that guided learning employs references without complete dependence on them. Conversely, residual learning exhibits failures or severe deviations when predicated on references of inferior quality. This shortfall stems from the framework's dependency on adhering closely to the provided references, which narrows its learning capabilities.

Nonetheless, reference-based learning reliance on predefined trajectories confines the policy to specific gaits, restricting its capacity to explore a broader range of motion possibilities. Additionally, this approach exhibits limited adaptability in responding effectively to unforeseen environmental changes or novel challenges.

## B. Reference-free learning

In reference-free learning, the policy is trained using a carefully crafted reward function rather than relying on predefined trajectories. This approach allows the policy to explore a wider range of gait patterns and adapt to unforeseen terrains, thereby enhancing innovation and flexibility within the learning process.

The concept of reference-free learning was initially explored using simulated physics engines with somewhat unrealistic bipedal models. A pioneering framework, which focused on learning symmetric gaits from scratch without the use of motion capture data, was developed and validated within a simulated environment [14]. This framework introduced a novel term into the loss function and utilized a curriculum learning strategy to effectively shape gait patterns. Another significant advancement was made in developing a learning method that enabled a robot to navigate stepping stones using curriculum learning, focusing on a physical robot model, Cassie, though this has yet to be validated outside of simulation [40].

Considering the practical implementation of theoretical models, significant efforts have been directed towards developing sim-to-real frameworks in robotics studies. A notable example of such a framework accommodates various periodic motions, including walking, hopping, and galloping [19]. This framework employs periodic rewards to facilitate initial training within simulations before successfully transitioning to a physical robot. It has been further refined to adapt to diverse terrains and scenarios. For instance, robust blind walking on stairs was demonstrated through terrain randomization techniques in [38]. Additionally, the integration of a vision system has enhanced the framework's ability to precisely determine foot locations [41], thus enabling the robot to effectively navigate stepping stones [20]. Subsequent developments include the incorporation of a vision system equipped with height maps, leading to an end-to-end framework that more effectively generalizes terrain information [42].

This approach to learning enables the exploration of novel solutions and strategies that might not be achievable through mere imitation of existing behaviours. However, the absence of reference guidance can render the learning process costly, time-consuming, and potentially infeasible for certain tasks. Moreover, the success of this method hinges critically on the design of the reward function, which presents significant challenges in specifying tasks such as jumping.

## III. HIERARCHY FRAMEWORK

Unlike end-to-end policies that directly map sensor inputs to motor outputs, hierarchical control schemes deconstruct locomotion challenges into discrete, manageable layers or stages of decision-making. Each layer within this structure is tasked with specific objectives, ranging from high-level navigation to fundamental locomotion skills. This division not only enhances the framework's flexibility but also simplifies the problem-solving process for each policy.

The architecture of a hierarchical framework typically comprises two principal modules: an HL planner and an LL controller. This modular approach allows for the substitution of each component with either a model-based method or a learning-based policy, further enhancing adaptability and customisation to specific needs.

Hierarchical frameworks can be classified into three distinct types based on the integration and function of their components:

1) **Deep planning hybrid scheme:** This approach combines strategic, high-level planning with dynamic low-level execution, leveraging the strengths of both learning-based and traditional model-based methods.
2) **Feedback DRL control hybrid scheme:** It focuses on integrating direct feedback control mechanisms with deep reinforcement learning, allowing for real-time adjustments and enhanced responsiveness.
3) **Learned hierarchy scheme:** Entirely learning-driven, this scheme develops a layered decision-making hierarchy where each level is trained to optimise specific aspects of locomotion.

These frameworks are illustrated in Fig. 3. Each type offers unique capabilities and exhibits distinct characteristics, albeit with limitations primarily due to the complexities involved in integrating diverse modules and their interactions.

For a concise overview, Table 3 summarises the various frameworks, detailing their respective strengths, limitations, and primary characteristics. The subsequent sections will delve deeper into each of these frameworks, providing a thorough analysis of their operational mechanics and their application in real-world scenarios.

## A. Deep planning hybrid scheme

In this scheme, robots are pre-equipped with the ability to execute basic locomotion skills such as walking, typically managed through model-based feedback controllers or interpretable methods. The addition of an HL learned layer focuses on strategic goals or the task space, enhancing locomotion capabilities and equipping the robot with advanced navigation abilities to effectively explore its environment.

Several studies have demonstrated the integration of an HL planner policy with a model-based controller to achieve tasks in world space. A notable framework optimises task space level performance, eschewing direct joint level and balancing considerations [24]. This system combines a residual learning planner with an inverse dynamics controller, enabling precise control over task-space commands to joint-level actions, thereby improving velocity tracking, foot touchdown location, and height control. Further advancements include a hybrid framework that merges HZD-based residual deep planning with model-based regulators to correct errors in learned trajectories, showcasing robustness, training efficiency, and effective velocity tracking [25]. These frameworks have been successfully transferred from simulation to reality and validated on robots such as Cassie.

However, the limitations imposed by residual learning constrained the agents' capacity to explore a broader array of possibilities. Building on previous work [25], a more efficient hybrid framework was developed, which learns from scratch

TABLE II: Summary and comparison of Hierarchy framework

| Control Scheme | Works | Module | characteristic | Advantages and Disadvantages |
|---|---|---|---|---|
| Deep Planning Hybrid Scheme | [24] [43] [44] | Deep planning + ID Deep planning + ID-QP Deep planning + WPG | HL policy is learned to guide the LL controller to complete locomotion and navigation tasks. | **A:** Enhanced command tracking capabilities, generalized across different platforms, sampling efficiency, and robust. **D:** Complicate system and communication between layers, require precise model, lack generalization regarding different tasks. |
| Feedback DRL Control Hybrid Scheme | [45] [26] [27] | Gait library + feedback policy Footstep planner + feedback policy Model-based planner + feedback policy | LL feedback policy receives non-learned HL planner as input to achieve locomotion skills. | **A:** Short inference times, robust, navigation locomotion capabilities, interpretability. **D:** Complicated system and communication between layers, reducing sampling efficiency. |
| Learned Hierarchy Framework | [12] [13] [29] | HL policy + LL policy HL policy + LL policy HL policy + LL policy | Both HL planner and LL feedback controller are learned. LL policy focuses on basic locomotion skills; on the other side, HL policy learn navigation skills. | **A:** provides layer flexibility, where each layer can be independently retrained and reused; alleviates the challenges associated with training an end-to-end policy. **D:** inefficiently sim-to-real, complicated interface between layers, training expensively. |

without reliance on prior knowledge [43]. In this approach, a purely learning-based HL planner interacts with an LL controller using an Inverse Dynamics with Quadratic Programming formulation (ID-QP). This policy adeptly captures dynamic walking gaits through the use of reduced-order states and simplifies the learning trajectory. Demonstrating robustness and training efficiency, this framework has outperformed other models and was successfully generalized across various bipedal platforms, including Digit, Cassie, and RABBIT.

In parallel, several research teams have focused on developing navigation planners specifically for toy-like humanoid robots, which provide greater physical stability compared to torque-driven or hydraulic bipedal robots as shown in Fig. 1. One notable study [46] implemented a visual navigation policy on the NAO robot, depicted in Fig. 1(a), utilizing RGB cameras as the primary sensory modality. This system has demonstrated successful zero-shot transfer to real-world scenarios, enabling the robot to adeptly navigate around obstacles. Further research [44] has explored complex dynamic motion tasks, such as playing soccer, by integrating a learned policy with an online footstep planner that utilises weight positioning generation (WPG) to create a center of mass (CoM) trajectory. This configuration is coupled with a whole-body controller, facilitating dynamic activities like soccer shooting. Despite their platform's stability, provided by large feet and a lightweight structure, these robots exhibit limited dynamic movement capabilities compared to full-sized humanoid robots. Consequently, this research primarily addresses navigation and task execution.

Regarding generalization, these frameworks have shown potential for adaptation across different types of bipedal and humanoid robots with minimal adjustments, demonstrating advanced user command tracking [43] and sophisticated navigation capabilities [44]. However, limitations are evident, notably the absence of capabilities for executing more complex

and dynamic motions, such as jumping. Furthermore, while these systems adeptly navigate complex terrains with obstacles, footstep planning alone is insufficient without concurrent enhancements to the robot's overall locomotion capabilities. Moreover, the requisite communication between the two distinct layers of the hierarchical framework may introduce system complexities. Enhancing both navigation and dynamic locomotion capabilities within the HL planner remains a significant challenge.

### B. Feedback DRL control hybrid scheme

In contrast to the comprehensive approach of end-to-end policies discussed in Section II, which excels in handling versatile locomotion skills and complex terrains with minimal interface times, the Feedback DRL Control Hybrid Scheme integrates DRL policies as LL controllers. These LL controllers, replacing traditional model-based feedback mechanisms, work in conjunction with HL planners that process terrain information, plan future walking paths, and maintain robust locomotion stability.

For instance, gait libraries, which provide predefined movement references based on user commands, have been integrated into such frameworks [45]. Despite the structured approach of using gait libraries, their static nature offers limited adaptability to changing terrains, diminishing their effectiveness. A more dynamic approach involves online planning, which has shown greater adaptability and efficiency. One notable framework combines a conventional foot planner with an LL DRL policy [26], delivering targeted footsteps and directional guidance to the robot, thereby enabling responsive and varied walking commands. Moreover, HL controllers can provide additional feedback to LL policies, incorporating CoM or end-feet information, either from model-based methods or other conventional control strategies. However, this work has not yet been transferred from simulation to real-world applications.
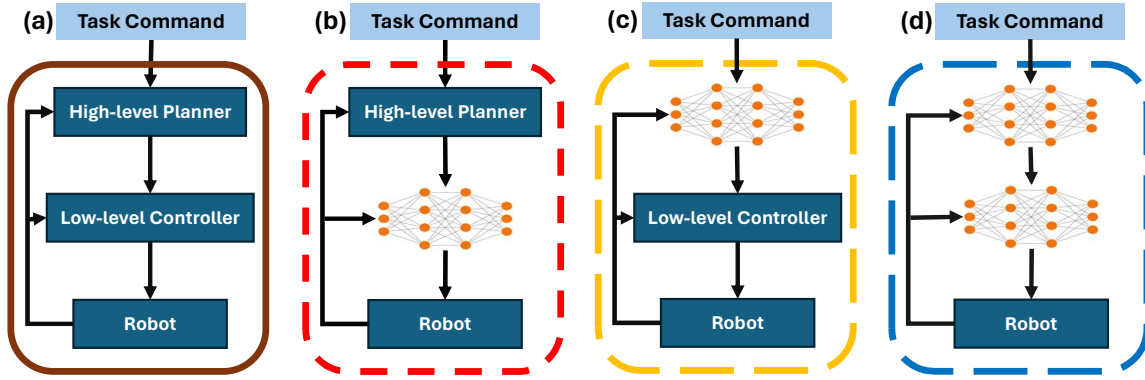
Fig. 3: **Hierarchy Control Scheme Diagram:** (a) A basic hierarchical scheme with two layers, where each module can be substituted with a learned policy. (b) A deep planning hybrid scheme, where the High-Level (HL) planner is learned. (c) A learning-based feedback control hybrid scheme, with a learned Low-Level (LL) controller. (d) A comprehensive DRL hierarchy control scheme, where both layers are learned.

Later, a similar structure featuring an HL foot planner and an LL DRL policy was proposed [27]. This strategy not only achieved a successful sim-to-real transfer but also enabled the robot to navigate omnidirectionally and avoid obstacles.

A recent development has shown that focusing solely on foot placement might restrict the stability and adaptability of locomotion, particularly in complex maneuvers. A new framework integrates a model-based planner with a DRL feedback policy to enhance bipedal locomotion's agility and versatility, displaying improved performance [47]. This system employs a residual learning architecture, where the DRL policy's outputs are merged with the planner's directives before being relayed to the PD controller. This integrated approach not only concerns itself with foot placement but also generates comprehensive trajectories for trunk position, orientation, and ankle yaw angle, enabling the robot to perform a wide array of locomotion skills including walking, squatting, turning, and stair climbing.

Compared to traditional model-based controllers, learned DRL policies provide a comprehensive closed-loop control strategy that does not rely on assumptions about terrain or robotic capabilities. These policies have demonstrated high efficiency in locomotion and accurate reference tracking. Despite their extensive capabilities, such policies generally require short inference times, making DRL a preferred approach in scenarios where robustness is paramount or computational resources on the robot are limited. Nonetheless, these learning algorithms often face challenges in environments characterized by sparse rewards, where suitable footholds like gaps or stepping stones are infrequent [48].

Additionally, an HL planner can process critical data such as terrain variations or obstacles and generate precise target locations for feet or desired walking paths, instead of detailed terrain data, which can significantly expedite the training process [27]. This capability effectively addresses the navigational limitations observed in end-to-end frameworks. Moreover, unlike the deep planning hybrid scheme where modifications post-policy establishment can be cumbersome, this hybrid scheme offers enhanced flexibility for on-the-fly adjustments.

Despite the significant potential demonstrated by previous studies, integrating DRL-based controllers with sophisticated and complex HL planners still presents limitations compared to more integrated frameworks such as end-to-end and deep planning models. Specifically, complex HL model-based planners often require substantial computational resources to resolve problems, rely heavily on model assumptions, necessitate extensive training periods, demand large datasets for optimization and hinder rapid deployment and iterative enhancements [48].

### C. Learned hierarchy framework

The Learned Hierarchy Framework merges a learned HL planner with an LL controller, focusing initially on refining LL policies to ensure balance and basic locomotion capabilities. Subsequently, an HL policy is developed to direct the robot towards specific targets, encapsulating a structured approach to robotic autonomy.

The genesis of this framework was within a physics engine, aimed at validating its efficiency through simulation [12]. In this setup, LL policies, informed by human motions or trajectories generated via Trajectory Optimization (TO), strive to track these trajectories as dictated by the HL planner while maintaining balance. An HL policy is then introduced, pre-trained with long-term task goals, to navigate the environment and identify optimal paths. This structure enabled sophisticated interactions such as guiding a biped to dribble a soccer ball towards a goal. The framework was later enhanced to include imitation learning, facilitating the replication of dynamic human-like movements within the simulation environment [13].

However, despite its structured and layered approach, which allows for the reuse of learned behaviors to achieve long-term objectives, these frameworks have predominantly been validated only in simulations. The interface designed manually between the HL planner and the LL controller sometimes leads to suboptimal behaviors, including stability issues like falling.

Expanding the application of this framework, a sim-to-real strategy for a wheeled bipedal robot was proposed, focusing the LL policy on balance and position tracking, while the HL policy enhances safety by aiding in collision avoidance and making strategic decisions based on the orientation of subgoals [29].

Learning complex locomotion skills, particularly when incorporating navigation elements, presents a significant challenge in robotics. Decomposing these tasks into distinct locomotion and navigation components allows robots to tackle more intricate activities, such as dribbling a soccer ball [12]. As discussed in the previous section, the benefits of integrating RL-based planners with RL-based controllers have been effectively demonstrated. This combination enables the framework to adeptly manage a diverse array of environments and tasks.

Within such a framework, the High-Level (HL) policy is optimized for strategic planning and achieving specific goals. This optimization allows for targeted enhancements depending on the tasks at hand. Moreover, the potential for continuous improvement and adaptation through further training ensures that the system can evolve over time, improving its efficiency and effectiveness in response to changing conditions or new objectives.

Despite the theoretical advantages, the practical implementation of this type of sim-to-real application for bipedal robots remains largely unexplored. The transition from simulation to real-world scenarios is fraught with challenges, not least because of the complexities involved in training and integrating two separate layers within the control hierarchy. Ensuring effective communication and cooperation between these layers is critical, requiring a meticulously defined communication interface to avoid operational discrepancies.

Additionally, the training process for each policy within the hierarchy demands considerable computational resources. The intensive nature of this training can lead to a reliance on the simulation environment, potentially causing the system to overfit to specific scenarios and thereby fail to generalize to real-world conditions. This limitation highlights a significant hurdle that must be addressed to enhance the viability of learned hierarchy frameworks in practical applications.

## IV. CHALLENGES AND FUTURE RESEARCH DIRECTIONS

While learning-based frameworks for bipedal robots have demonstrated considerable potential, they have also clearly exposed the limitations inherent to each framework. Moreover, several critical areas remain largely unexplored, especially within the realm of legged robotics, where the pace of research on bipedal robots lags behind that of their quadruped counterparts. This discrepancy in research progress can be attributed to several factors, including the higher costs and less mature technology associated with bipedal robot hardware, as well as the inherent instability issues that bipedal designs face.

To gain a deeper understanding of these challenges and to outline potential future directions, it is instructive to first review existing research on quadruped robots. The insights gained from quadrupeds, which benefit from more robust research outputs and technological advancements, can provide

valuable lessons for overcoming similar challenges in bipedal systems.

### A. Recent progress with quadruped robots

While DRL remains an emerging technology in bipedal robotics, it has firmly established its presence in the realm of quadruped robots, another category of legged systems. The diversity of frameworks developed for quadrupeds ranges from model-based RL designed for training in real-world scenarios, where unpredictable dynamics often prevail [49], [50], to systems that include the modeling of deformable terrain to enhance locomotion over compliant surfaces [51]. Furthermore, dynamic quadruped models facilitate highly adaptable policies [52], and sophisticated acrobatic motions are achieved through imitation learning [53].

The domain of quadruped DRL has also seen significant advancements in complex hybrid frameworks that integrate vision-based systems. To date, two primary versions of such frameworks have been developed: one where a deep planning module is paired with model-based control [54], and another that combines model-based planning with low-level DRL control [48], [55]. The latter has shown substantial efficacy; it employs a model predictive control (MPC) to generate reference motions, which are then followed by a LL feedback DRL policy. Additionally, the Terrain-aware Motion Generation for Legged Robots (TAMOLS) module [56] enhances the MPC and DRL policy by providing terrain height maps for effective foothold placements across diverse environments, including those not encountered during training. However, similar hybrid control schemes have not been thoroughly investigated within the field of bipedal locomotion.

Quadruped DRL frameworks are predominantly designed to navigate complex terrains, but efforts to extend their capabilities to other tasks are underway. These include mimicking real animals through motion capture data and imitation learning [57], [58], as well as augmenting quadrupeds with manipulation abilities. This is achieved either by adding a manipulator [59], [60] or by using the robots' legs [61]. Notably, the research presented in [60] demonstrates that loco-manipulation tasks can be effectively managed using a single unified end-to-end framework.

Despite the progress in quadruped DRL, similar advancements have been limited for bipedal robots, particularly in loco-manipulation tasks and vision-based DRL frameworks. Establishing a unified framework could bridge this gap, an essential step given the integral role of bipedal robots in developing full humanoid systems. Moreover, the potential of hybrid frameworks that combine model-based and DRL-based methods in bipedal robots remains largely untapped.

### B. Gaps and challenges

Despite numerous promising developments in the field of bipedal and humanoid robotics, significant gaps remain between current research outcomes and the ultimate goals. This discussion concentrates on the gaps in frameworks and algorithms rather than hardware, structured around two pivotal questions: 1) Is it possible to design a unified framework that

achieves both generalization and precision? 2) Can we develop a straightforward end-to-end policy capable of managing all tasks efficiently?

*1) Generalization versus precision:* DRL has demonstrated potential in facilitating versatile locomotion skills [22]; however, challenges such as poor velocity tracking and issues with precise control often arise. While [43] shows that deep planning combined with model-based control can achieve precise velocity tracking, and [37] illustrates successful end-to-end control for precise jumping, the creation of a policy that effectively handles both diverse tasks and precise movements remains elusive. Furthermore, [41] introduces a foot constraint policy framework, enabling precise target tracking and accurate touchdown locations. Yet, there is still no framework that comprehensively addresses the dual demands of versatility and precision in locomotion.

The difficulty in simultaneously achieving precise control and a broad range of actions in bipedal locomotion using DRL stems from several factors:

- **Complex dynamics:** Bipedal locomotion involves intricate dynamics, posing a significant challenge to maintaining both dynamic motion and precision.
- **Resource intensity:** Executing diverse locomotion tasks requires considerable computational power and extensive data, necessitating high-quality hardware and efficient DRL algorithms.
- **Training conflicts:** Training DRL systems to achieve both precision and versatility often leads to conflicts. Designing reward functions and training policies that satisfy both criteria is inherently complex.

These challenges underscore the need for innovative solutions that can bridge the gap between the capabilities of current frameworks and the ambitious goals of advanced bipedal and humanoid robotics.

*2) Simplifying frameworks to overcome complex tasks:* The envisioned ideal in robotic design is an end-to-end framework that enables robots to traverse various terrains using versatile locomotion skills. Although current research often focuses on enhancing frameworks by adding complex components to mitigate inherent limitations, such as the integration of a foot planner for omnidirectional locomotion and stair navigation, as demonstrated in [26], [27], simpler end-to-end frameworks have also proven effective. These frameworks adeptly navigate challenging terrains and perform a diverse range of locomotion tasks with fewer components [20], [22].

The advantage of maintaining simplicity in the framework lies in its ability to streamline decision-making processes, thereby reducing computational overhead and potential points of failure. To achieve an optimal end-to-end framework, advancements in several key areas are essential:

- **Robust and efficient DRL algorithms:** Development of algorithms that can manage high-dimensional and continuous control problems more effectively.
- **Specialized neural network architectures:** Design of neural architectures tailored for specific bipedal tasks, capable of processing extensive sensory data (e.g., visual and tactile inputs), similar to the innovations presented in [42].

- **Effective reward functions:** Formulation of reward functions that more accurately guide the learning process towards achieving desired behaviors and strategic outcomes.
- **Advanced computational resources:** Enhancement of computational capabilities to support more intensive training and faster inference, facilitating real-time decision-making in dynamic environments.

By focusing on these developmental areas, the potential to create a unified, efficient, and less complex framework for handling complex locomotion challenges in bipedal robots is significantly increased.

## C. Future directions

The exploration of quadruped robotics has yielded substantial advancements, yet the full potential of bipedal robotics remains largely untapped. Building on the successes and innovative approaches observed in quadruped robots, several key future directions emerge that could significantly enhance bipedal and humanoid robotics.

*1) Unified framework:* Currently, no single framework exists that enables bipedal or humanoid robots to adeptly navigate all types of terrains, including stepping stones, stairs, deformable terrain, and slippery surfaces. A promising approach, as evidenced by recent work in quadruped robots [48], utilizes MPC to generate reference motions, which a low-level DRL policy then tracks. This method, coupled with the Terrain-aware Motion Generation for Legged Robots (TAMOLS) module, simplifies the terrain representation into a height map, facilitating more effective navigation. This success encourages further exploration into hybrid frameworks that combine model-based methods with DRL, inheriting the strengths of both approaches, as discussed in Section III. However, hybrid frameworks present challenges such as training efficiency and system complexity, which demand considerable computational resources and extensive training periods.

Moreover, recent studies [20], [42], [22] have demonstrated the potential of end-to-end frameworks enhanced with vision-based information. These frameworks successfully navigate challenging terrains and execute dynamic motions, suggesting the feasibility of a unified framework capable of handling diverse environments and tasks. Training strategies such as curriculum learning and task randomization could be employed, utilizing visual height maps as inputs to the policy, enhancing the robot's ability to adapt and perform in varied scenarios.

In addition, the introduction of a DRL end-to-end framework incorporating transformer models, as in [62], presents significant possibilities for integrating locomotion skills with language and vision capabilities. The use of large-scale models capable of processing and condensing extensive data sets into a coherent model could expand the robot's range of capabilities, maintaining versatility across a broad spectrum of tasks.

The exploration of transformers and other large-scale models holds considerable promise for enhancing generalizability and adaptability in complex tasks, warranting further investigation into their potential applications in bipedal robotics.

*2) Vision-based learning framework:* Vision plays a critical role in enabling robots to navigate challenging terrains, such as blind drops, where tactile and other sensory inputs may not provide sufficient information. Despite the importance of vision, many current frameworks, particularly in bipedal robotics, do not fully exploit this modality [38], [43]. Vision-based systems are essential in human locomotion for identifying obstacles and assessing terrains, and some studies have begun to show the effectiveness of integrating vision into DRL frameworks for bipedal and humanoid robots [20], [42], [27].

Building on the groundwork laid by both bipedal and quadruped robots, two promising directions have emerged:

- **Height scanner mapping:** This approach, evaluated in works like [42], involves using height maps generated by scanners to inform locomotion strategies. These maps provide detailed topographical data, allowing robots to plan steps on uneven or obstructed surfaces more effectively.
- **Direct vision inputs:** Directly utilizing inputs from cameras, such as depth or RGB images, for real-time decision-making in RL policies [46], [63]. Although previous studies like [46] have integrated visual navigation by feeding visual information to a High-Level (HL) planner, the potential of direct visual inputs to RL policies has not been fully explored.

Enhancing the capability of bipedal robots to directly interpret and utilize visual data without intermediary processing can revolutionize their adaptability and efficiency in real-world scenarios. The exploration of direct vision inputs to reinforcement learning policies represents a significant opportunity for advancing the field, potentially enabling more dynamic and responsive locomotion strategies.

*3) Bridge the gap from simulation to reality:* While simulations offer a safe and cost-effective environment for developing robotics policies, the transition from simulation to real-world application often encounters significant challenges due to the approximations and simplifications made in simulations. Numerous sim-to-real frameworks [34], [64], [65] have shown high efficiency and performance, as detailed in Appendix B. Despite these advancements, a significant gap persists, exacerbated by the complexity and unpredictability of physical environments. Moreover, many studies [26], [66], [21] remain validated only in simulation settings.

*4) Loco-manipulation tasks:* Loco-manipulation, which combines locomotion and manipulation, presents opportunities for humanoid robots to excel beyond purely bipedal capabilities. Few studies have addressed this integrated task; one such study [67] demonstrated a 'box transportation' framework. This framework decomposes the task into five distinct policies, each addressing different aspects of the transportation process. However, this approach lacks efficiency and does not incorporate vision-based information, suggesting substantial room for improvement. Moreover, the challenges of managing mobile tools like scooters [68] or dynamically interacting with objects such as balls [69] introduce further complexities.

Decomposing loco-manipulation tasks into multiple layers could simplify the challenges, allowing for more precise and flexible control by manually tuning individual components of the task [43]. This structured control approach provides a more coordinated response to complex interactions within the robot's environment, facilitating the execution of task-specific commands.

Alternatively, an end-to-end framework may enable bipeds to perform a variety of tasks through task randomization and structured curriculum learning methods, progressively teaching the policy [35], [27], [22]. During training, such policies can also learn human-like movements from motion capture data [70], [18], [16], offering promising solutions for future integrated loco-manipulation tasks within a single, versatile policy.

*5) Designing reward functions:* The development of effective reward functions is a critical challenge in the field of deep reinforcement learning (DRL) for bipedal robots. While periodic reward functions have been designed to facilitate cyclic movements like walking [19], there remains a significant gap in crafting reward functions for non-periodic actions such as jumping. These actions require distinct considerations for success and efficiency, yet current research lacks comprehensive methods for their reward structure. Furthermore, minimizing the need for extensive manual tuning while achieving high performance in DRL systems continues to be a substantial challenge, pointing to the need for more adaptive and automatically adjusting reward mechanisms.

*6) Integrating large language models:* The integration of Large Language Models (LLMs) into bipedal robotics opens new avenues for contextual understanding and task execution, significantly enhancing the robots' interaction capabilities. LLMs, when implemented at the highest task level, offer substantial promise for improving human-robot interaction, making these systems more intuitive and responsive [71]. The potential applications of this technology are broad and impactful, spanning sectors such as industrial automation, where robots can perform complex assembly tasks; healthcare, offering assistance in patient care and rehabilitation; assistive devices, providing support for individuals with disabilities; search and rescue operations, where robust and adaptive decision-making is critical; and entertainment and education, where interactive and engaging experiences are key [72]. Each of these fields could benefit from the advanced capabilities of LLM-enhanced bipedal robots, particularly in environments requiring nuanced understanding and adaptability.

### D. Applications in various fields

The advancements in bipedal locomotion technology hold significant promise for practical applications beyond the confines of laboratory environments. These robots, bolstered by AI, are poised to transform numerous sectors by enhancing operational capabilities and interaction with humans. The potential for humanoid robots in various fields is detailed in [72], emphasizing the integration of learning-based approaches for more effective implementation. Key areas include:

1) **Industrial automation and manufacturing**: The integration of humanoid robots in industrial settings can significantly enhance productivity and efficiency, freeing workers from repetitive and labor-intensive tasks. These

robots, equipped with advanced loco-manipulation capabilities and the ability to cooperate with human teams, are particularly effective in assembly line operations, maintenance tasks, and the construction of complex machinery [73], [74]. Their articulated arms and floating bases provide unmatched flexibility, making them ideal for human-centric manufacturing environments. The humanoid robot Digit, for example, demonstrates remarkable stability and efficiency in industrial tasks over extended periods, as seen in video demonstrations [75]. Moreover, these robots are also suited for operation in high-risk environments such as underwater or areas with high radiation levels, significantly enhancing safety and operational capacity in these contexts.

2) **Healthcare and assistive devices**: In the healthcare sector, bipedal and humanoid robots contribute significantly to rehabilitation and assistive technologies. Exoskeletons enhanced with DRL methodologies are being used to train individuals to achieve more natural gait patterns, improving mobility and rehabilitation outcomes [76]. Beyond mere mobility aids, humanoid robots integrated with LLMs show promise in delivering medications, monitoring patient health, and assisting in surgeries [77]. The synergy between LLMs and loco-manipulation capabilities paves the way for more interactive, responsive support, aligning closely with the needs of personalized care. Additionally, the aging population can benefit from humanoid robots performing everyday tasks like house cleaning or delivery through simple voice commands, thereby enhancing the quality of life.

3) **Search and rescue missions**: Humanoid robots are exceptionally valuable in search and rescue operations, especially in disaster-stricken or hazardous environments where human presence is risky or impractical. Unlike traditional wheeled robots, humanoid robots can navigate complex terrains filled with debris, gaps, and elevated structures, making them indispensable in these scenarios. They also demonstrate potential for significant interaction and collaboration with human rescue teams. For instance, in environments with high nuclear radiation, humanoid robots can perform tasks that would be perilous for humans, handling delicate instruments and preventing human exposure to harmful conditions. This capability extends to other challenging environments such as underwater [78], outer space [79], [80], and other hazardous areas. However, the full realization of these applications remains constrained by the absence of a unified framework that can seamlessly navigate all terrains and fully integrate loco-manipulation and human interaction functionalities.

4) **Entertainment and education**: Humanoid robots have the potential to transform the realms of entertainment and education by providing highly interactive experiences. With their ability to integrate extensive knowledge bases, these robots can significantly enhance educational environments. They can assume the roles of butlers, teachers, or even babysitters, engaging with users in diverse activities. For example, robots can facilitate language learning [81], participate in storytelling, teach various academic subjects, or engage in the performing arts and games. In the sphere of entertainment, humanoid robots can act, dance, play ball games [69], and take part in interactive performances, captivating audiences of all ages with their versatility and dynamic capabilities.

However, this in turn leads to a variety of ethical issues. First, interacting with humans involves collecting human daily behavior data and increases the risk of a data breach. Second, another concern is the increasing dependency of humans on robots, not just for assistance but also for emotional support. This will result in less human-to-human interaction and ultimately affect social constructs and emotional development. Third, advancements of humanoid robots will replace humans in various jobs, and eventually lead to unemployment issues.

On the positive side, humanoid robots can provide invaluable assistance to people with disabilities or the elderly, offering companionship and reducing the care burden on families and healthcare systems. Furthermore, their application across diverse fields such as education, industry, and healthcare can bring about revolutionary changes, improving efficiency and safety while opening up new possibilities for technological integration. As we navigate these advancements, it is crucial to balance innovation with ethical considerations to ensure that the deployment of humanoid robots enhances societal well-being without compromising personal integrity or social dynamics.

## V. Conclusion

Despite significant advances in DRL for robotics, a considerable gap persists between current research achievements and the development of a unified framework capable of empowering robots to perform a broad spectrum of complex tasks efficiently. Presently, DRL research can be categorized into two primary control schemes: end-to-end and hierarchical frameworks. End-to-end frameworks have shown promising capabilities in executing diverse locomotion skills [22], climbing stairs [38], and navigating challenging terrains such as stepping stones [20]. Conversely, hybrid frameworks, which often integrate an HL planner or an LL model-based controller, offer enhanced capabilities, allowing for simultaneous management of locomotion and navigation tasks.

To bridge the existing gaps, further development of hierarchical frameworks, particularly those equipped with advanced perception systems and integrated with model-based planners, appears promising. Such frameworks could simultaneously address issues of precision and generalization. Moreover, the advent of LLMs presents a transformative opportunity, potentially enabling the unification of language processing and visual functionalities within robotic systems. While numerous challenges remain—ranging from the technical intricacies of framework integration to real-world application—the steady progression in control framework refinement and DRL development provides a hopeful outlook. The vision of achieving an

end-to-end unified framework, capable of mimicking human-like learning processes and enabling bipedal robots to handle a wide range of complex tasks, may soon move within reach.

## APPENDIX A
### DEEP REINFORCEMENT LEARNING ALGORITHMS

The advancement and development of RL is crucial for bipedal locomotion. Specifically, advancements in deep learning provide deep neural networks serving as function approximators to empower RL with the capability to handle tasks characterized by high-dimensional and continuous spaces, by efficiently discovering condensed, low-dimensional representations of complex data. In comparison to other robots of different morphologies, such as wheeled robots, bipedal robots feature much higher DoFs and continuously interact with environments, which results in higher requirements for the DRL algorithms. Especially in the legged locomotion field, policy gradient-based algorithms are prevalent in the field of bipedal locomotion.

Designing an effective neural network architecture is essential for tackling complex bipedal locomotion tasks. Multilayer perceptrons (MLP), a fundamental neural network, excel in straightforward regression tasks with lower computational resource demands. A comprehensive comparison between MLP and the memory-based neural network, Long Short-Term Memory (LSTM) reveals that MLPs have an advantage in convergence speed for tasks [65]. However, LSTM, as a variant of Recurrent Neural Networks (RNN), is adept at processing data associated with time, effectively relating different states across time, and modeling key physical properties vital for periodical gaits [19] and successful sim-to-real transfer in bipedal locomotion. Additionally, Convolutional Neural Networks (CNN) specialize in spatial data processing, particularly for image-related tasks, making them highly suitable for environments where visual perception is crucial. This diverse range of neural network architectures highlights the importance of selecting the appropriate model based on the specific requirements of the bipedal locomotion tasks.

Considering DRL alogrithms, recent bipedal locomotion studies focus on model-free reinforcement algorithms. Unlike model-based RL, which learns a model of the environment but may inherit biases from simulations that do not accurately reflect real-world conditions, model-free RL directly trains policies through environmental interaction without relying on an explicit environmental model. Although model-free RL requires more computational samples and resources, it can train a more robust policy allowing the robots to travel around challenging environments.

Many sophisticated model-free RL algorithms exist, which can be broadly classified into two categories: policy-based (or policy optimization) and value-based approaches. Value-based methods e.g. Q-learning, Deep Q-learning (DQN) [82] only excel in discrete action space and often struggle with high dimensional action space. In contrast, policy-based methods, such as policy gradient, can handle complex tasks but are generally less sample-efficient compared to value-based methods.

More advanced algorithms combine both policy-based methods and value-based methods. Actor-critic (AC) refers
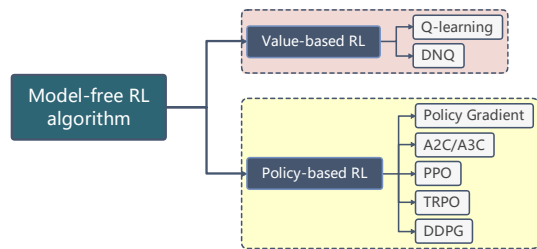


Fig. 4: **Diagram for RL algorithms catalogue**

to a main idea simultaneously learning both a policy (actor) and a value function (critic), where it owns both advantages of both algorithms [83], [84]. Popular algorithms e.g. Trust region policy optimization (TRPO) [85] and PPO based on policy-based methods, borrow ideas from AC. Moreover, there are other novel algorithms based on the AC framework, Deep Deterministic Policy Gradient (DDPG) [86], Twin Delayed Deep Deterministic Policy Gradients (TD3) [87], A2C (Advantage Actor-Critic), and A3C (Asynchronous Advantage Actor-Critic) [88], SAC (Soft Actor-Critic) [89]. Each algorithm has its strengths considering different tasks in the bipedal locomotion scenario. There are several key factors to value these algorithms such as: sample efficiency, robustness and generalization, and implementation challenges. A comparative analysis work [90] illustrates that SAC-based algorithms excel in stability and achieve the highest scores, while their training efficiency significantly trails behind that of PPO that obtain relatively high score.

In [91], PPO demonstrates the robustness and computational economy in complex scenarios, such as bipedal locomotion, utilizing fewer resources than TRPO. In terms of training time, PPO is much faster than SAC, and DDPG algorithms [90]. Besides, many works [19], [45], [36] have demonstrated its robustness and ease of implementation and combined with the flexibility to integrate with various neural network architectures have made PPO the most popular choice in this field. Various work has demonstrated that PPO can conduct the exploration of walking [19], jumping [37], stair climbing [38], and stepping stones [20], which demonstrates its efficiency, robustness and generalization.

Additionally, the DDPG algorithm integrates the Actor-Critic framework with DQN to facilitate off-policy training, further optimizing sampling efficiency. In some explicit scenarios such as jumping, DDPG shows higher reward and better learning performance than PPO [21], [92]. TD3 is developed based on DDPG, and improve over the performance of the DDPG and SAC [89]. Soft Actor-Critic (SAC) further the agent's exploration capabilities and sample efficiency [89]. While A2C offers improved efficiency and stability compared to A3C, the asynchronous update mechanism for A3C provides better capabilities for exploration and accelerating learning. Although these algorithms show their advancements, they are more challenging to apply due to algorithms' complexity compared to PPO.

## APPENDIX B
### BRIDGING SIM-TO-REAL GAP

Due to the large number of interactions needed for RL algorithms, training directly on robots can lead to costly damage to hardware and the environment. Consequently, training a policy in the simulation and then deploying it to the hardware illustrates significant potential and efficiency. However, the gap between simulation and the real world remains substantial, making sim-to-real challenging. To overcome the gap, several sim-to-real approaches are developed, including dynamics randomization [36], system identification [93], [94], periodic reward composition [66], learned actuators dynamics [93], [95], regulation feedback controller [96], adversarial motion prior [18], [17].

There are two primary approaches to training policies under domain randomization. One is end-to-end training with a history of robot measurements or I/O [36] and another is policy distillation, an expert policy with environmental insights guides a student policy that learns from internal sensory feedback, such as teacher-student policy [66], RMA [64].

Details of these sim-to-real transition approaches are shown below:

1) The dynamics randomization method involves systematically varying the physical parameters of the simulated environment-such as mass, inertia, or stiffness. 2) System Identification methods develop mathematical models of dynamics from observed data, enhancing the accuracy of robots' properties within the models, such as mass and inertia, to ensure the model faithfully represents the system's behavior. 3) The learned actuator dynamics method utilizes experimental data from actuators to develop a model of their dynamics, achieving sim-to-real by incorporating realistic actuator behavior within the training environment. It is noticeable that the higher-level planner can also learn from reference or non-reference. 4) periodic reward composition helps capture the essential locomotion information and the periodic gait pattern is more general to adapt to uncertainty and variation in the real world. 5) The regulation feedback controller manually tunes the setting of the controller to mitigate perturbations and gaps between the sim and the real, thereby enhancing the robustness and adaptation. Key aspects of sim-to-real, including system identification, state estimation with noise measurement, and the selection of state-and-action spaces are highlighted in [34].

### REFERENCES

[1] S. Gupta and A. Kumar, "A brief review of dynamics and control of underactuated biped robots," *Advanced Robotics*, vol. 31, pp. 607–623, 2017.

[2] J. Reher and A. Ames, "Dynamic walking: Toward agile and efficient bipedal robots," *Annual Review of Control, Robotics, and Autonomous Systems*, vol. 4, 2021.

[3] J. Carpentier and P.-B. Wieber, "Recent progress in legged robots locomotion control," *Current Robotics Reports*, vol. 2, pp. 231–238, 2021.

[4] M. A.-M. Khan, M. R. J. Khan, A. Tooshil, N. Sikder, M. A. P. Mahmud, A. Z. Kouzani, and A.-A. Nahid, "A systematic review on reinforcement learning-based robotics within the last decade," *IEEE Access*, vol. 8, pp. 176 598–176 623, 2020.

[5] J. García and D. Shafie, "Teaching a humanoid robot to walk faster through safe reinforcement learning," *Engineering Applications of Artificial Intelligence*, vol. 88, p. 103360, 2020.

[6] C. Chevallereau, G. Abba, Y. Aoustin, F. Plestan, E. Westervelt, C. C. De Wit, and J. Grizzle, "Rabbit: A testbed for advanced control theory," *IEEE Control Systems Magazine*, vol. 23, pp. 57–79, 2003.

[7] Y. Gong, R. Hartley, X. Da, A. Hereid, O. Harib, J.-K. Huang, and J. Grizzle, "Feedback control of a cassie bipedal robot: Walking, standing, and riding a segway," in *American Control Conference*, 2019, pp. 4559–4566.

[8] S. Kuindersma, R. Deits, M. Fallon, A. Valenzuela, H. Dai, F. Permenter, T. Koolen, P. Marion, and R. Tedrake, "Optimization-based locomotion planning, estimation, and control design for the atlas humanoid robot," *Autonomous robots*, vol. 40, pp. 429–455, 2016.

[9] G. A. Castillo, B. Weng, W. Zhang, and A. Hereid, "Robust feedback motion policy design using reinforcement learning on a 3d digit bipedal robot," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2021, pp. 5136–5143.

[10] R. Tedrake, T. Zhang, and H. Seung, "Stochastic policy gradient reinforcement learning on a simple 3D biped," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2004, pp. 2849–2854.

[11] J. Morimoto, G. Cheng, C. Atkeson, and G. Zeglin, "A simple reinforcement learning algorithm for biped walking," in *IEEE International Conference on Robotics and Automation*, 2004, pp. 3030–3035 Vol.3.

[12] X. Peng, G. Berseth, K. Yin, and M. Panne, "DeepLoco: dynamic locomotion skills using hierarchical deep reinforcement learning," *ACM Transactions on Graphics*, vol. 36, pp. 1–13, 2017.

[13] X. Peng, P. Abbeel, S. Levine, and M. Panne, "DeepMimic: Example-guided deep reinforcement learning of physics-based character skills," *ACM Transactions on Graphics*, vol. 37, 2018.

[14] W. Yu, G. Turk, and C. K. Liu, "Learning symmetric and low-energy locomotion," *ACM Transactions on Graphics*, vol. 37, pp. 1–12, 2018.

[15] M. Taylor, S. Bashkirov, J. F. Rico, I. Toriyama, N. Miyada, H. Yanagisawa, and K. Ishizuka, "Learning bipedal robot locomotion from human movement," in *IEEE International Conference on Robotics and Automation*, 2021, pp. 2797–2803.

[16] X. Cheng, Y. Ji, J. Chen, R. Yang, G. Yang, and X. Wang, "Expressive whole-body control for humanoid robots," *arXiv preprint arXiv:2402.16796*, 2024.

[17] A. Tang, T. Hiraoka, N. Hiraoka, F. Shi, K. Kawaharazuka, K. Kojima, K. Okada, and M. Inaba, "HumanMimic: Learning natural locomotion and transitions for humanoid robot via wasserstein adversarial imitation," *arXiv preprint arXiv:2309.14225*, 2023.

[18] Q. Zhang, P. Cui, D. Yan, J. Sun, Y. Duan, A. Zhang, and R. Xu, "Whole-body humanoid robot locomotion with human reference," *arXiv preprint arXiv:2402.18294*, 2024.

[19] J. Siekmann, Y. Godse, A. Fern, and J. Hurst, "Sim-to-real learning of all common bipedal gaits via periodic reward composition," in *IEEE International Conference on Robotics and Automation*, 2021, pp. 7309–7315.

[20] H. Duan, A. Malik, M. S. Gadde, J. Dao, A. Fern, and J. Hurst, "Learning dynamic bipedal walking across stepping stones," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2022, pp. 6746–6752.

[21] C. Tao, M. Li, F. Cao, Z. Gao, and Z. Zhang, "A multiobjective collaborative deep reinforcement learning algorithm for jumping optimization of bipedal robot," *Advanced Intelligent Systems*, vol. 6, p. 2300352, 2023.

[22] Z. Li, X. B. Peng, P. Abbeel, S. Levine, G. Berseth, and K. Sreenath, "Reinforcement learning for versatile, dynamic, and robust bipedal locomotion control," *arXiv e-prints*, pp. arXiv–2401, 2024.

[23] T. Li, H. Geyer, C. G. Atkeson, and A. Rai, "Using deep reinforcement learning to learn high-level policies on the ATRIAS biped," in *International Conference on Robotics and Automation*, 2019, pp. 263–269.

[24] H. Duan, J. Dao, K. Green, T. Apgar, A. Fern, and J. Hurst, "Learning task space actions for bipedal locomotion," in *IEEE International Conference on Robotics and Automation*, 2021, pp. 1276–1282.

[25] G. A. Castillo, B. Weng, W. Zhang, and A. Hereid, "Reinforcement learning-based cascade motion policy design for robust 3d bipedal locomotion," *IEEE Access*, vol. 10, pp. 20 135–20 148, 2022.

[26] R. P. Singh, M. Benallegue, M. Morisawa, R. Cisneros, and F. Kanehiro, "Learning bipedal walking on planned footsteps for humanoid robots," in *IEEE-RAS International Conference on Humanoid Robots*, 2022, pp. 686–693.

[27] S. Wang, S. Piao, X. Leng, and Z. He, "Learning 3D bipedal walking with planned footsteps and fourier series periodic gait planning," *Sensors*, vol. 23, p. 1873, 2023.

[28] K. Arulkumaran, M. P. Deisenroth, M. Brundage, and A. A. Bharath, "Deep reinforcement learning: A brief survey," *IEEE Signal Processing Magazine*, vol. 34, pp. 26–38, 2017.

[29] W. Zhu and M. Hayashibe, "A hierarchical deep reinforcement learning framework with high efficiency and generalization for fast and safe navigation," *IEEE Transactions on Industrial Electronics*, vol. 70, pp. 4962–4971, 2023.

[30] X. B. Peng and M. Van De Panne, "Learning locomotion skills using deeprl: Does the choice of action space matter?" in *ACM SIG-GRAPH/Eurographics Symposium on Computer Animation*, 2017, pp. 1–13.

[31] Z. Li, X. Cheng, X. B. Peng, P. Abbeel, S. Levine, G. Berseth, and K. Sreenath, "Reinforcement learning for robust parameterized locomotion control of bipedal robots," in *IEEE International Conference on Robotics and Automation*, 2021, pp. 2811–2817.

[32] D. Kim, G. Berseth, M. Schwartz, and J. Park, "Torque-based deep reinforcement learning for task-and-robot agnostic learning on bipedal robots using sim-to-real transfer," *IEEE Robotics and Automation Letters*, vol. 8, p. 6251–6258, 2023.

[33] Z. Xie, G. Berseth, P. Clary, J. Hurst, and M. van de Panne, "Feedback control for cassie with deep reinforcement learning," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2018, pp. 1241–1246.

[34] Z. Xie, P. Clary, J. Dao, P. Morais, J. Hurst, and M. van de Panne, "Learning locomotion skills for cassie: Iterative design and sim-to-real," in *Conference on Robot Learning*, 2020, pp. 317–329.

[35] D. Rodriguez and S. Behnke, "Deepwalk: Omnidirectional bipedal gait by deep reinforcement learning," in *IEEE International Conference on Robotics and Automation*, 2021, pp. 3033–3039.

[36] J. Siekmann, S. Valluri, J. Dao, L. Bermillo, H. Duan, A. Fern, and J. W. Hurst, "Learning memory-based control for human-scale bipedal locomotion," in *Robotics science and systems*, 2020.

[37] Z. Li, X. B. Peng, P. Abbeel, S. Levine, G. Berseth, and K. Sreenath, "Robust and versatile bipedal jumping control through multi-task reinforcement learning," in *Robotics: Science and Systems*, 2023.

[38] J. Siekmann, K. Green, J. Warila, A. Fern, and J. Hurst, "Blind bipedal stair traversal via sim-to-real reinforcement learning," in *Robotics: Science and Systems*, 2021.

[39] 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*, 2019, pp. 270–276.

[40] Z. Xie, H. Ling, N. Kim, and M. Panne, "ALLSTEPS: Curriculum-driven learning of stepping stone skills," *Computer Graphics Forum*, vol. 39, pp. 213–224, 2020.

[41] H. Duan, A. Malik, J. Dao, A. Saxena, K. Green, J. Siekmann, A. Fern, and J. Hurst, "Sim-to-real learning of footstep-constrained bipedal dynamic walking," in *International Conference on Robotics and Automation*, 2022, pp. 10 428–10 434.

[42] B. Marum, M. Sabatelli, and H. Kasaei, "Learning vision-based bipedal locomotion for challenging terrain," *arXiv preprint arXiv:2309.14594*, 2023.

[43] G. A. Castillo, B. Weng, S. Yang, W. Zhang, and A. Hereid, "Template model inspired task space learning for robust bipedal locomotion," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2023, pp. 8582–8589.

[44] C. Gaspard, G. Passault, M. Daniel, and O. Ly, "FootstepNet: an efficient actor-critic method for fast on-line bipedal footstep planning and forecasting," *arXiv preprint arXiv:2403.12589*, 2024.

[45] K. Green, Y. Godse, J. Dao, R. L. Hatton, A. Fern, and J. Hurst, "Learning spring mass locomotion: Guiding policies with a reduced-order model," *IEEE Robotics and Automation Letters*, vol. 6, pp. 3926–3932, 2021.

[46] K. Lobos-Tsunekawa, F. Leiva, and J. Ruiz-del Solar, "Visual navigation for biped humanoid robots using deep reinforcement learning," *IEEE Robotics and Automation Letters*, vol. 3, no. 4, pp. 3247–3254, 2018.

[47] J. Li, L. Ye, Y. Cheng, H. Liu, and B. Liang, "Agile and versatile bipedal robot tracking control through reinforcement learning," *arXiv preprint arXiv:2404.08246*, 2024.

[48] F. Jenelten, J. He, F. Farshidian, and M. Hutter, "DTC: Deep tracking control," *Science Robotics*, vol. 9, p. eadh5401, 2024.

[49] L. Smith, I. Kostrikov, and S. Levine, "Demonstrating a walk in the park: Learning to walk in 20 minutes with model-free reinforcement learning," *Robotics: Science and Systems Demo*, vol. 2, p. 4, 2023.

[50] P. Wu, A. Escontrela, D. Hafner, P. Abbeel, and K. Goldberg, "Day-Dreamer: World models for physical robot learning," in *Conference on Robot Learning*, 2023, pp. 2226–2240.

[51] 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, p. eade2256, 2023.

[52] G. Feng, H. Zhang, Z. Li, X. B. Peng, B. Basireddy, L. Yue, Z. SONG, L. Yang, Y. Liu, K. Sreenath, and S. Levine, "Genloco: Generalized locomotion controllers for quadrupedal robots," in *Conference on Robot Learning*, vol. 205, 2023, pp. 1893–1903.

[53] Y. Fuchioka, Z. Xie, and M. Van de panne, "OPT-Mimic: Imitation of optimized trajectories for dynamic quadruped behaviors," in *IEEE International Conference on Robotics and Automation*, 2023, pp. 5092–5098.

[54] S. Gangapurwala, M. Geisert, R. Orsolino, M. Fallon, and I. Havoutis, "RLOC: Terrain-aware legged locomotion using reinforcement learning and optimal control," *IEEE Transactions on Robotics*, vol. 38, pp. 2908–2927, 2022.

[55] D. Kang, J. Cheng, M. Zamora, F. Zargarbashi, and S. Coros, "RL + Model-Based Control: Using on-demand optimal control to learn versatile legged locomotion," *IEEE Robotics and Automation Letters*, vol. 8, pp. 6619–6626, 2023.

[56] F. Jenelten, R. Grandia, F. Farshidian, and M. Hutter, "TAMOLS: Terrain-aware motion optimization for legged systems," *IEEE Transactions on Robotics*, vol. 38, pp. 3395–3413, 2022.

[57] X. B. Peng, E. Coumans, T. Zhang, T.-W. E. Lee, J. Tan, and S. Levine, "Learning agile robotic locomotion skills by imitating animals," in *Robotics: Science and Systems*, 2020.

[58] F. Yin, A. Tang, L. Xu, Y. Cao, Y. Zheng, Z. Zhang, and X. Chen, "Run like a dog: Learning based whole-body control framework for quadruped gait style transfer," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2021, pp. 8508–8514.

[59] Y. Ma, F. Farshidian, T. Miki, J. Lee, and M. Hutter, "Combining learning-based locomotion policy with model-based manipulation for legged mobile manipulators," *IEEE Robotics and Automation Letters*, vol. 7, pp. 2377–2384, 2022.

[60] Z. Fu, X. Cheng, and D. Pathak, "Deep whole-body control: Learning a unified policy for manipulation and locomotion," in *Conference on Robot Learning*, 2023, pp. 138–149.

[61] P. Arm, M. Mittal, H. Kolvenbach, and M. Hutter, "Pedipulate: Enabling manipulation skills using a quadruped robot's leg," in *IEEE Conference on Robotics and Automation*, 2024.

[62] I. Radosavovic, T. Xiao, B. Zhang, T. Darrell, J. Malik, and K. Sreenath, "Learning humanoid locomotion with transformers," *arXiv preprint arXiv:2303.03381*, 2023.

[63] A. Byravan, J. Humplik, L. Hasenclever, A. Brussee, F. Nori, T. Haarnoja, B. Moran, S. Bohez, F. Sadeghi, B. Vujatovic *et al.*, "Nerf2real: Sim2real transfer of vision-guided bipedal motion skills using neural radiance fields," in *IEEE International Conference on Robotics and Automation*, 2023, pp. 9362–9369.

[64] A. Kumar, Z. Li, J. Zeng, D. Pathak, K. Sreenath, and J. Malik, "Adapting rapid motor adaptation for bipedal robots," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2022, pp. 1161–1168.

[65] R. P. singh, Z. Xie, P. Gergondet, and F. Kanehiro, "Learning bipedal walking for humanoids with current feedback," *IEEE Access*, vol. 11, p. 82013–82023, 2023.

[66] B. van Marum, M. Sabatelli, and H. Kasaei, "Learning perceptive bipedal locomotion over irregular terrain," *arXiv preprint arXiv:2304.07236*, 2023.

[67] J. Dao, H. Duan, and A. Fern, "Sim-to-real learning for humanoid box loco-manipulation," *arXiv preprint arXiv:2310.03191*, 2023.

[68] J. Baltes, G. Christmann, and S. Saeedvand, "A deep reinforcement learning algorithm to control a two-wheeled scooter with a humanoid robot," *Engineering Applications of Artificial Intelligence*, vol. 126, p. 106941, 2023.

[69] T. Haarnoja, B. Moran, G. Lever, S. H. Huang, D. Tirumala, J. Humplik, M. Wulfmeier, S. Tunyasuvunakool, N. Y. Siegel, R. Hafner *et al.*, "Learning agile soccer skills for a bipedal robot with deep reinforcement learning," *Science Robotics*, vol. 9, p. eadi8022, 2024.

[70] M. Seo, S. Han, K. Sim, S. H. Bang, C. Gonzalez, L. Sentis, and Y. Zhu, "Deep imitation learning for humanoid loco-manipulation through human teleoperation," in *IEEE-RAS International Conference on Humanoid Robots*, 2023, pp. 1–8.

[71] K. N. Kumar, I. Essa, and S. Ha, "Words into action: Learning diverse humanoid robot behaviors using language guided iterative motion refinement," in *Workshop on Language and Robot Learning: Language as Grounding*, 2023.

[72] Y. Tong, H. Liu, and Z. Zhang, "Advancements in humanoid robots: A comprehensive review and future prospects," *IEEE/CAA Journal of Automatica Sinica*, vol. 11, pp. 301–328, 2024.

[73] A. Dzedzickis, J. Subačiūtė-Žemaitienė, E. Šutinys, U. Samukaitė-Bubnienė, and V. Bučinskas, "Advanced applications of industrial

robotics: New trends and possibilities," *Applied Sciences*, vol. 12, p. 135, 2021.

[74] M. Yang, E. Yang, R. C. Zante, M. Post, and X. Liu, "Collaborative mobile industrial manipulator: a review of system architecture and applications," in *International conference on automation and computing*, 2019, pp. 1–6.

[75] "6+ Hours Live Autonomous Robot Demo," https://www.youtube.com/watch?v=Ke468Mv8ldM, Mar. 2024.

[76] G. Bingjing, H. Jianhai, L. Xiangpan, and Y. Lin, "Human–robot interactive control based on reinforcement learning for gait rehabilitation training robot," *International Journal of Advanced Robotic Systems*, vol. 16, p. 1729881419839584, 2019.

[77] A. Diodato, M. Brancadoro, G. De Rossi, H. Abidi, D. Dall'Alba, R. Muradore, G. Ciuti, P. Fiorini, A. Menciassi, and M. Cianchetti, "Soft robotic manipulator for improving dexterity in minimally invasive surgery," *Surgical innovation*, vol. 25, pp. 69–76, 2018.

[78] R. Bogue, "Underwater robots: a review of technologies and applications," *Industrial Robot: An International Journal*, vol. 42, pp. 186–191, 2015.

[79] N. Rudin, H. Kolvenbach, V. Tsounis, and M. Hutter, "Cat-like jumping and landing of legged robots in low gravity using deep reinforcement learning," *IEEE Transactions on Robotics*, vol. 38, pp. 317–328, 2022.

[80] J. Qi, H. Gao, H. Su, L. Han, B. Su, M. Huo, H. Yu, and Z. Deng, "Reinforcement learning-based stable jump control method for asteroid-exploration quadruped robots," *Aerospace Science and Technology*, vol. 142, p. 108689, 2023.

[81] O. Mubin, C. Bartneck, L. Feijs, H. Hooft van Huysduynen, J. Hu, and J. Muelver, "Improving speech recognition with the robot interaction language," *Disruptive science and Technology*, vol. 1, pp. 79–88, 2012.

[82] A. Meduri, M. Khadiv, and L. Righetti, "DeepQ stepper: A framework for reactive dynamic walking on uneven terrain," in *IEEE International Conference on Robotics and Automation*, 2021, pp. 2099–2105.

[83] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," in *International Conference on Learning Representations*, 2016.

[84] L. Liu, M. V. D. Panne, and K. Yin, "Guided learning of control graphs for physics-based characters," *ACM Transactions on Graphics*, vol. 35, pp. 1–14, 2016.

[85] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz, "Trust region policy optimization," in *International Conference on Machine Learning*, 2015, pp. 1889–1897.

[86] C. Huang, G. Wang, Z. Zhou, R. Zhang, and L. Lin, "Reward-adaptive reinforcement learning: Dynamic policy gradient optimization for bipedal locomotion," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, pp. 7686–7695, 2023.

[87] S. Dankwa and W. Zheng, "Twin-delayed DDPG: A deep reinforcement learning technique to model a continuous movement of an intelligent robot agent," in *International conference on vision, image and signal processing*, 2019, pp. 1–5.

[88] J. Leng, S. Fan, J. Tang, H. Mou, J. Xue, and Q. Li, "M-A3C: A mean-asynchronous advantage actor-critic reinforcement learning method for real-time gait planning of biped robot," *IEEE Access*, vol. 10, pp. 76 523–76 536, 2022.

[89] C. Yu and A. Rosendo, "Multi-modal legged locomotion framework with automated residual reinforcement learning," *IEEE Robotics and Automation Letters*, vol. 7, pp. 10 312–10 319, 2022.

[90] O. Aydogmus and M. Yilmaz, "Comparative analysis of reinforcement learning algorithms for bipedal robot locomotion," *IEEE Access*, pp. 7490–7499, 2023.

[91] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv e-prints*, pp. arXiv–1707, 2017.

[92] C. Tao, J. Xue, Z. Zhang, and Z. Gao, "Parallel deep reinforcement learning method for gait control of biped robot," *IEEE Transactions on Circuits and Systems II: Express Briefs*, vol. 69, pp. 2802–2806, 2022.

[93] W. Yu, V. C. V. Kumar, G. Turk, and C. K. Liu, "Sim-to-real transfer for biped locomotion," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2019, pp. 3503–3510.

[94] S. Masuda and K. Takahashi, "Sim-to-real transfer of compliant bipedal locomotion on torque sensor-less gear-driven humanoid," in *IEEE-RAS International Conference on Humanoid Robots*, 2023, pp. 1–8.

[95] 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, p. eaau5872, 2019.

[96] G. A. Castillo, B. Weng, W. Zhang, and A. Hereid, "Robust feedback motion policy design using reinforcement learning on a 3D digit bipedal robot," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2021, pp. 5136–5143.