# Adaptive Control Strategy for Quadruped Robots in Actuator Degradation Scenarios

Xinyuan Wu*
Shanghai Jiao Tong University
Shanghai, China
wuxinyuan@sjtu.edu.cn

Wentao Dong*
Shanghai Jiao Tong University
Shanghai, China
a_Dong@sjtu.edu.cn

Hang Lai
Shanghai Jiao Tong University
Shanghai, China
laihang99@sjtu.edu.cn

Yong Yu
Shanghai Jiao Tong University
Shanghai, China
yyu@apex.sjtu.edu.cn

Ying Wen†
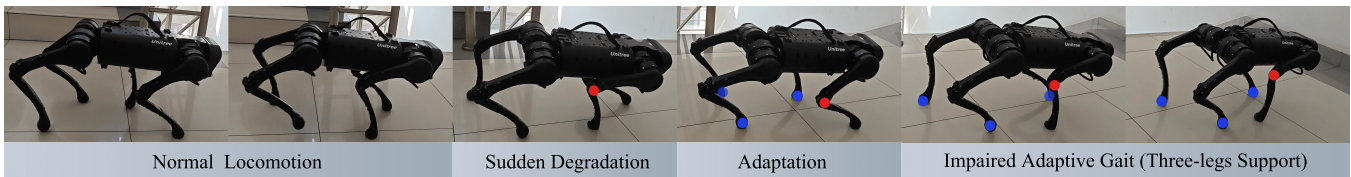Shanghai Jiao Tong University
Shanghai, China
Ying.Wen@sjtu.edu.cn

Figure 1: Application of ADAPT facing sudden degradation (indicated by red).

## ABSTRACT

Quadruped robots have strong adaptability to extreme environments but may also experience faults. Once these faults occur, robots must be repaired before returning to the task, reducing their practical feasibility. One prevalent concern among these faults is actuator degradation, stemming from factors like device aging or unexpected operational events. Traditionally, addressing this problem has relied heavily on intricate fault-tolerant design, which demands deep domain expertise from developers and lacks generalizability. Learning-based approaches offer effective ways to mitigate these limitations, but a research gap exists in effectively deploying such methods on real-world quadruped robots. This paper introduces a pioneering teacher-student framework rooted in reinforcement learning, named **A**ctuator **D**egradation **A**da**p**tation **T**ransformer (ADAPT), aimed at addressing this research gap. This framework produces a unified control strategy, enabling the robot to sustain its locomotion and perform tasks despite sudden joint actuator faults, relying exclusively on its internal sensors. Empirical evaluations on the Unitree A1 platform validate the deployability and effectiveness of ADAPT on real-world quadruped robots, and affirm the robustness and practicality of our approach.

---

*Both authors contributed equally to this research.

†corresponding author.

---

## KEYWORDS

Deep Reinforcement Learning, Quadruped Robots, Machine Learning for Robot Control, Fault Tolerance, Real-World Deployment.

## 1 INTRODUCTION

Legged robots, known for their multifunctionality, have demonstrated remarkable adaptability and flexibility in traversing complex terrains and navigating diverse and unfamiliar environments [3, 27]. As hardware technology and control algorithms advance, quadruped robots have evolved to possess improved load-bearing capacities, heightened stability, and structural advantages [5]. These advancements have resulted in their growing deployment across diverse domains such as exploration [27], search and rescue [20], military applications [17], and industrial scenarios [3].

However, in real-world applications, the robustness of quadruped robots is critically compromised by their susceptibility to various fault scenarios. While these robots excel in navigating complex and unpredictable environments [5], their adaptability renders them vulnerable to various fault scenarios, such as physical damage, joint malfunctions, and motor aging. Such vulnerabilities can have far-reaching implications, leading to actuator degradation and consequent reductions in operational efficiency [18, 24]. Moreover, these faults not only inconvenience operators but, in extreme cases, may also jeopardize the longevity of the robot and pose significant safety risks to humans. Therefore, how to deal with such a fault has drawn significant attention from researchers.

Dealing with the fault in real-time is exceptionally challenging due to obstacles like distance-related accessibility issues and hazardous conditions. Additionally, immediate fault rectification often necessitates specialized expertise and tools, which may not be readily available. As a result, the design of fault-tolerant control algorithms becomes crucial, empowering robots to adapt to varying degrees of malfunction autonomously [10, 11]. Implementing such algorithms not only prolongs the robot's operational lifespan but also substantially enhances its reliability. This, in turn, broadens the range of real-world applications where quadruped robots can be effectively deployed.

While traditional fault-tolerant algorithms like Model Predictive Control (MPC) and whole-body control frameworks have been widely adopted, they suffer from concrete drawbacks such as heavy reliance on domain expertise, labor-intensive manual tuning, and limited adaptability to unforeseen environmental and mechanical conditions [11, 12, 18, 24]. These conventional approaches focus on the active detection and confirmation of faults during the robot's operation. If a fault is identified, a pre-engineered alternative control scheme is invoked to mitigate the issue. However, the effectiveness of these strategies is intrinsically tied to the designer's comprehensive understanding of the robot's mechanical intricacies, often necessitating meticulous manual adjustments [15]. Moreover, the rigidity of relying on predefined alternative control schemes becomes a significant constraint in real-world applications. Specifically, these methods struggle to adapt to novel environmental or mechanical scenarios, frequently requiring the formulation of separate, customized control strategies for different variants of the same fault [24].

While learning-based approaches, particularly DRL, offer substantial advantages such as autonomous skill acquisition and reduced dependency on domain expertise, they also present specific challenges, including sample efficiency, algorithmic stability, and the complexity of bridging the simulation-to-reality gap [2, 22, 26].

This paradigm empowers robots to acquire emergency-response skills independently [15], reducing the dependence on domain knowledge when designing algorithms to some extent [1, 24]. Deep Reinforcement Learning (DRL) has gained prominence in quadruped robot control due to its ability to learn nonlinear control strategies [26]. Unlike traditional control methods, DRL does not require precise models and effectively manages high-dimensional state spaces [2]. To balance exploration and exploitation, reinforcement learning algorithms consistently develop strategies with enhanced adaptability and generalization capabilities [22, 26, 43, 45]. Furthermore, the availability of physically accurate modeling simulation environments like IsaacGym [25] and massively parallel solutions [36] has significantly elevated the efficiency of control strategy learning while easing the transfer from simulation to real-world environments.

While substantial progress has been made in enhancing robot motion performance and enabling cross-terrain capabilities under stable operational conditions [21, 23], there remains a notable gap in exploring actuator degradation fault within the context of DRL for real-world quadruped robots.

In response to bridge the gap, we present a pioneering teacher-student training framework that combines transformer-based architectures with reinforcement learning principles, named **A**ctuator **D**egeneration **A**da**p**tation **T**ransformer (ADAPT). This framework equips the quadruped robot to adapt to joint actuator degradation faults within predetermined thresholds, relying solely on its perceptual abilities. Our proposed framework demonstrates remarkable robustness within simulated environments and showcases the potential for zero-shot transfer to real-world robots. This learning-based strategy reduces the dependence on deep expertise in the field of robotics or specific robot models while also holding significant potential for acquiring a broader range of skills.

This paper contributes three key aspects to quadruped robot control amidst joint actuator degradation faults.

- Firstly, we tackle an issue of varying degrees of actuator degradation in an intricate simulation environment. We approach this as a type of cross-embodiment task using reinforcement learning techniques that can handle various actuator degradation scenarios instead of certain pre-defined cases.
- Secondly, we introduce ADAPT, a novel adaptation framework that demonstrates significant generalization potential in experiments.
- Lastly, we successfully deployed ADAPT on real-world quadruped robots and impressively demonstrated the capability for zero-shot transferring from simulation to real-world robots.

These contributions collectively pave the way for more resilient and versatile quadruped robot locomotion strategies. To the best of our knowledge, this is the first work to successfully deploy a learning-based approach in addressing actuator degradation on real-world quadruped robots.

## 2 RELATED WORK

### 2.1 DRL for Robot Control

Deep Reinforcement Learning (DRL) holds the potential to reduce the dependence on expertise in robot control by training a policy in simulation and then transferring it to real world. However, it's challenging to directly apply the DRL policy to real world due to the discrepancy between simulation and real world, also known as *reality gap* [6, 19]. A lot of works were devoted to bridge the reality gap. For example, Domain Randomization (DR) [30, 40] proposes to train policy in a wide range of environments with various parameters and noises, which makes policy more robust. Additionally, one can also adapt the policy trained in simulator to real-world data in latent space, called domain adaptation [14, 39]. Besides, system identification, which tries to identify the physical parameters of real-world explicit or implicit [31, 46], may also help improving the transferring performance.

By leveraging these advanced skills, DRL has achieved remarkable strides in the fields of legged robot locomotion [21, 22, 31], robotic arm manipulation [4, 35, 38], and wheels robot application [8, 47, 48]. Within the domain of quadruped robots, Lee et al. [23] and Kumar et al. [21] developed the teacher-student training paradigm for quadruped robot learning, which has been demonstrated to be highly effective in experiments. Lai et al. [22] introduced a two-stage training framework called TERT, demonstrating its superior performance across diverse, challenging terrains, showcasing robust task mobility. Nahrendra et al. [28] presented a robust framework for quadrupedal locomotion, enabling stable movement across unstructured terrains. Peng et al. [31] presented a robust

framework allowing legged robots to acquire agile locomotion skills by imitating real-world animals' movements. Escontrela et al. [13] presented an effective style reward introducing Adversarial Motion Prior approach into quadrupted robot control. Similarly, Wu et al. [43] formulated a single policy, trained through DRL, to achieve a harmonious balance between robust and agile quadrupedal locomotion. Meanwhile, Yu et al. [45] developed a learning-based control method to control quadrupedal robots with varying morphologies effectively.

Notably, the works mentioned above primarily revolve around the premise of the robot's optimal operation, with seldom specific provisions for fault tolerance in the presence of malfunctions.

## 2.2 Fault-Tolerant Quadrupedal Locomotion

Common robot faults can be categorized into "locked" and "uncontrollable" [10]. In the former, the malfunctioning joint becomes immovable and can be used to provide support, while in the latter, the joint becomes freely movable and loses its supporting function. This paper mainly focuses on the "uncontrollable" fault scenario, which presents higher challenges and lacks a relevant learning-based framework for quadruped robots.

Numerous studies have delved into alternative strategies, including traditional control theories and other methodologies, to tackle the challenge of fault tolerance in quadruped robots. For instance, Cui et al. [11] and Chen et al. [10] introduced control-based methods to enable quadruped robots to overcome single-joint locking scenarios. Zhao et al. [49] devised an adaptive fault-tolerant control law to address leg joint actuator faults in quadruped robots. Allard et al. [1] and their team have undertaken a series of endeavors rooted in Quality-Diversity that successfully tackled the challenge of damage recovery in real-world hexapod robots. This approach, inherent in autonomous learning, showcases significant potential.

However, a common challenge these methods share is their substantial reliance on task-specific expertise, as they exhibit a heightened dependence on the specification of lower-level behavioral logic. In contrast, the end-to-end learning capabilities of deep reinforcement learning, coupled with the inherent adaptive autonomy, can yield a robust and stable control strategy to address robot faults. Okamoto et al. [29] demonstrated the feasibility of deep reinforcement learning in tackling actuator failure within the ant-v2 environment of the OpenAI Gym [7]. Yan et al. [44] introduced a fault-tolerant reinforcement learning control framework for robotic manipulator joint actuator faults. Furthermore, the DRL-based hardware fault-tolerant controller devised by Wu et al. [43] for quadruped robot locomotion successfully validates in both simulator and real-world environments. Importantly, their emphasis was on the "locked" fault scenario, as previously highlighted, which is notably distinct from the central focus of our investigation.

## 3 PRELIMINARIES

### 3.1 Reinforcement Learning for Robotic Locomotion

We conceptualize robotic locomotion as a Markov Decision Process (MDP), delineated by the tuple $(\mathcal{S}, \mathcal{A}, P, r, \gamma, \mu_0)$, where $\mathcal{S}$ denotes the robotic state space, containing the robot's proprioception information. The symbol $\mathcal{A}$ signifies the action space, which is described

in detail in Section 4.1. The transition density $P(\cdot|s_t, a_t)$ specifies the likelihood of transitioning to various states when taking action $a_t$ in the given state $s_t$. We denote the reward function as $r(s_t, a_t)$, while $\gamma$ represents the discount factor and $\mu_0$ represents the initial state distribution. Reinforcement Learning (RL) strives to identify the optimal policy $\pi^* : \mathcal{S} \rightarrow \mathcal{A}$, which maximizes the expected accumulated return, factoring in the discount, throughout decision-making:

$$\pi^* = \arg\max_{\pi} \mathbb{E}_{s_0 \sim \mu_0, a_t \sim \pi(\cdot|s_t), s_{t+1} \sim P(\cdot|s_t, a_t)} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right]. \quad (1)$$

### 3.2 Adversarial Motion Prior

The design of reward functions for quadrupedal locomotion is non-trivial and usually needs much domain knowledge and human effort [21, 23]. To address this issue, Adversarial Motion Prior (AMP) [13] is proposed to simplify the reward function design by utilizing some pre-collected expert motions as prior. Specifically, AMP utilizes a similar framework with generative adversarial imitation learning (GAIL) [16] to discriminate whether a transition $(s_t, s_{t+1})$ is sampled from expert motions. The training objective for the discriminator is defined as:

$$\arg\min_{\varphi} \mathbb{E}_{(s_t, s_{t+1}) \sim \mathcal{D}} \left[ (D_{\varphi}(s_t, s_{t+1} - 1)^2 \right]$$
$$+ \mathbb{E}_{(s_t, s_{t+1}) \sim \pi} \left[ (D_{\varphi}(s_t, s_{t+1}) + 1)^2 \right] \quad (2)$$
$$+ \frac{\alpha^{gp}}{2} \mathbb{E}_{(s_t, s_{t+1}) \sim \mathcal{D}} \left[ \|\nabla_{\varphi} D_{\varphi}(s_t, s_{t+1})\|_2 \right].$$

which contains two least square terms and a gradient penalty term with coefficient $\alpha^{gp}$. The discriminator is then used to calculate the style reward via Equation 3, which gives a larger reward if the state encountered has a lower probability of being discriminated from the expert dataset by $D_{\varphi}$.

$$r_t^{\text{style}}(s_t, s_{t+1}) = \max \left[ 0, 1 - 0.25(D_{\varphi}(s_t, s_{t+1}) - 1)^2 \right]. \quad (3)$$

Intuitively, The style function regulates agents' actions with a similar style of reference, which guides the robot into a natural gait.

### 3.3 Decision Transformer

The Transformer architecture [41] has garnered significant attention in various domains, including computer vision, natural language processing, and beyond. The attention mechanisms capture the long-range dependencies and relevance, effectively tackling sequential decision-making problems [42]. Recognizing its potential, Chen et al. abstracted the reinforcement learning (RL) problem into the framework of conditional sequence modeling, subsequently proposing a transformer-based solution known as Decision Transformer (DT) [9]. By leveraging the autoregressive model, DT predicts actions based on past state-action sequences and expected returns.

Consistent with prior research, we employ the GPT architecture [32], comprising several stacked self-attention layers with residual connections and causal masking. The input sequences are embedded as tokens containing position encoding. Subsequently, each token of this sequence is mapped to key, query, and value pairs $k_i$, $q_i$, and $v_i$ through the self-attention layers. The resultant output of the self-attention layer is represented as follows:

$$z_i = \sum_{j=1}^{n} \text{softmax}(\{\langle q_i, k_{j'} \rangle\}_{j'=1}^{n})_j \cdot v_j. \tag{4}$$

In addition, we use position encoding [41] to model the sequential relevance:

$$p_t^i = \begin{cases} \cos(\omega_k \cdot t) & \text{if } i = 2k \\ \sin(\omega_k \cdot t) & \text{o.w.} \end{cases}, \omega_k = \frac{1}{10000^{2k/d}}. \tag{5}$$
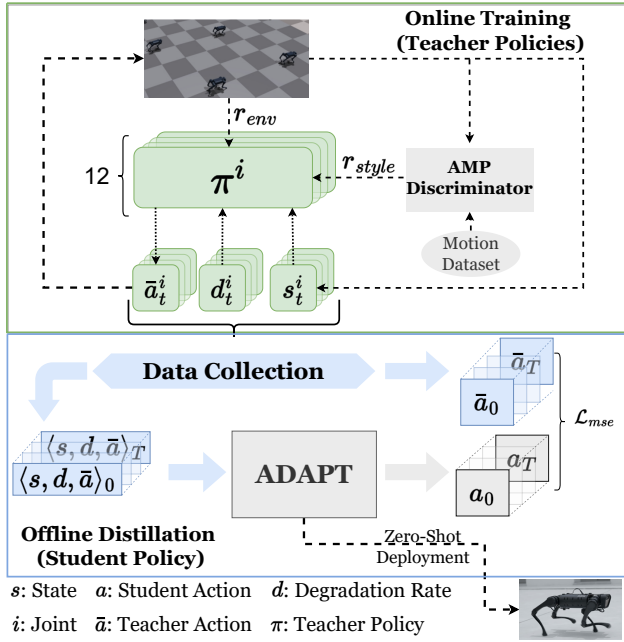
## 4 METHODOLOGY



**Figure 2: Overall framework of ADAPT. ADAPT starts by training 12 teacher policies separately in simulation. These teacher policies are then used to generate trajectories, which are subsequently utilized to distill a unified transformer-based student policy. Afterward, the student policy is poised for zero-shot deployment to real-world robots.**

This section comprehensively explains our novel approach, the Actuator Degradation Adaptation Transformer (ADAPT). ADAPT is devised as a teacher-student framework meticulously designed to incorporate and leverage actuator degradation information. The teacher policy within the ADAPT framework employs actuator degradation information as prior knowledge. This utilization of prior information serves the purpose of narrowing down the potential state space during policy exploration. On the other hand, the student policy within the ADAPT framework utilizes actuator degradation information to enhance decision-making processes.

### 4.1 RL-based Control Architecture

The environment is defined precisely as the Markov Decision Process (MDP) outlined in Section 3.1.

**State.** In the context of this article, our data originates solely from the robot's internal sensory inputs, without incorporating any external sources or additional information such as images or GPS data. Consistent with previous designs [22, 36, 45], we use a 48-dimensional vector $s \in \mathcal{S}$ as the state. This state vector contains both the base and each joint's information (12 joints in total). The base information includes several key variables:

- Base linear velocity $v \in \mathbb{R}^3$. Estimated through comprehensive readings from the IMU accelerometer, the current orientation of the base returned by IMU, and the velocity estimation of the contact leg.
- Base angular velocity $\omega \in \mathbb{R}^3$. Derived from the IMU gyroscope.
- The body's orientation, represented by the gravity component $g \in \mathbb{R}^3$, is calculated based on the base direction measured by the IMU.
- Target linear velocity commands $(\hat{v}_x, \hat{v}_y) \in \mathbb{R}^2$ and target angular velocity $\hat{\omega}_z \in \mathbb{R}$.

The information for each joint encompasses:

- Positions $q \in \mathbb{R}^{12}$ and velocities $\dot{q} \in \mathbb{R}^{12}$ can be retrieved through the motor encoders.
- The action $a' \in \mathcal{A}$ from the previous time step will also be considered.

**Action.** $a \in \mathcal{A}$, a 12-dimensional vector, signifies the target position of each joint. Based on insights from traditional control theory, position control is considered safer and more stable than direct torque control. We subsequently relay the action $a_t$ through a Proportional-Derivative (PD) controller, which generates the desired torque output for each actuator:

$$\tau_{\text{desired}}^i = K_p(\hat{q} - q) + K_d(\hat{\dot{q}} - \dot{q}). \tag{6}$$

Here, $(K_p, K_d)$ represents the stiffness and damping gains. These two hyper-parameters are fundamental components of the PD controller and remain constant throughout our entire training process.

**Actuator Degradation Rate.** Inspired by the design concept of EAT [45], we introduce the actuator degradation rate to quantify the degradation extent of quadruped robot joints at the current time step. Denoted as $d \in \mathcal{D} \subset \mathbb{R}^{12}$, this 12-dimensional vector represents the actuator degradation rate for each joint within the robot's structure. Notably, $d^i$ is continuous, effectively adapting to diverse fault scenarios during the quadruped robot's operation. This continuity enables smooth transfers between different fault states. Each component $d^i \in [0, 1]$, where $d^i = 0$ denotes an intact joint, and $d^i = 1$ signifies a joint entirely devoid of force. Formally, the actuator degradation can be expressed as:

$$\tau_{\text{applied}}^i = \tau_{\text{desired}}^i (1 - d^i). \tag{7}$$

Here, $\tau_{\text{applied}}^i$ represents the actual torque applied to joint $i$, while $\tau_{\text{desired}}^i$ is the desired torque on joint $i$. The PD controller determines the latter based on the specified action through Equation 7. Notably, this torque adjustment directly affects the actual torque output to the joint in both simulated and real-world settings.

**Reward Function.** Our reward function, denoted as $r : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$, follows the reward configuration outlined in the article [43]. However, in conjunction with experimental knowledge, our approach

**Table 1: Reward Function Settings**

| Term | Reward Formula | Scale |
|------|----------------|-------|
| Linear velocity tracking | $\exp\left(-\frac{(\mathbf{v}_{\text{cmd}} - \mathbf{v}_{\text{base}})^2}{0.25}\right)$ | 1.0 |
| Angular velocity tracking | $\exp\left(-\frac{(\omega_{\text{cmd}} - \omega_{\text{base}})^2}{0.25}\right)$ | 0.5 |
| Unexpected angular velocity | $\omega_{\text{pitch}}^2 + \omega_{\text{roll}}^2$ | -0.05 |
| Orientation | $x_{\text{gravity}}^2 + y_{\text{gravity}}^2$ | -2 |
| Torques | $\|\tau\|^2$ | -0.0001 |
| Accelerations | $(\mathbf{v}_{t-1} - \mathbf{v}_t)^2$ | -2.5e-7 |
| Feet air time | $\|\mathbf{t}_{\text{air}} - 0.5 \cdot \mathbf{1}\|_1$ | 1.0 |
| Collisions | $\mathbf{1}_{\text{unnatural touchdown}}$ | -0.5 |
| Changes in actions | $\left(a_{t-1}^n - a_t^n\right)^2$ | -0.01 |
| Large action | $\|a_t\|^2$ | -0.3 |

is further informed by insights from previous work [22, 36]. We have incorporated penalties targeting orientation deviations in the angular velocities of the base's roll and pitch. This strategic addition encourages the robot to adopt more stable gaits, particularly when navigating under fault conditions. Additionally, we introduce a reward term designed to incentivize the robot to take long footsteps. Concurrently, we impose penalties for robot collisions. These combined reward and penalty mechanisms contribute to our pursuit of refining the robot's behavior. Table 1 shows the reward function used to evaluate the agents.

## 4.2 Teacher Policy

Given the high sensitivity of the actuator degradation rate within the task, we train twelve distinct teacher conditional policies. Each policy corresponds to one of the twelve scenarios that one joint actuator degrades to some extent. This strategy effectively reduces the policy exploration space, enhancing efficiency. Moreover, we adopt the AMP architecture [13] as the foundation of our teacher policy, which regulates the gaits of quadrupedal robots, aligning them with the locomotion patterns observed in expert behaviors. Specifically, in addition to the reward terms outlined in Table 1, we introduce a term of style reward [13] (mentioned in Equation 3) that evaluates the similarity between the current strategy's gait and that of expert motion.

Let $\bar{\pi}^i$ denote the $i$-th teacher policy for the scenario where the $i$-th joint degrades. The actuator degradation rate at time step $t$ for the $i$-th joint is denoted as $d_t^i$. The value of $d_t^i$ is sampled from the continuous interval $[0, 1]$ at the beginning of each episode. Subsequently, the teacher policy $\bar{\pi}^i$ outputs actions conditioned on the actuator degradation rate $d_t^i$ and the corresponding state $s_t^i$:

$$\bar{a}_t^i = \bar{\pi}^i(s_t^i, d_t^i). \tag{8}$$

The teacher policies are trained via the Proximal Policy Optimization algorithm [37] and use domain randomization [40] to enhance real-world performance.

## 4.3 Data Collection for ADAPT Training

The meticulous selection of expert data proves pivotal within our teacher-student framework for knowledge distillation. Our experiment highlights a crucial insight: As the extent of degradation escalates, the robot's gait significantly transforms. Therefore, to facilitate the learning of the student, we sample $d_i$ uniformly from

$[0, 1]$ at timestep 0 and then adjust it adaptively during data collection:

$$d_t^i = \begin{cases} \mathcal{U}(d_{t-1}^i, 1) & \text{w.p. } p, \\ \mathcal{U}(0, 0.5) & d_{t-1}^i > 1 - \delta, \\ d_{t-1}^i & \text{o.w.} \end{cases} \tag{9}$$

where $\mathcal{U}$ denotes the uniform distribution, $\delta$ represents the reset threshold and $p$ signifies the probability of degradation. This mechanism ensures the collected data cover an extensive range of possible degradation scenarios, ensuring a broader skill acquisition by the students.

Using this mechanism, we collect a set of trajectories encompassing $12 \times N$ trajectories in total, each containing $T$ transitions:

$$\mathcal{T} = \{\mathcal{T}^i\}_{i=0,1,\dots,11},$$
$$\mathcal{T}^i = \{traj_0^i, traj_1^i, \dots, traj_N^i\}, \tag{10}$$
$$traj_k^i = \{(s_{k,t}^i, d_{k,t}^i, \bar{a}_{k,t}^i)\}_{t=0,1,\dots,T}.$$

Specifically, we set $N = 20000, T = 500, \delta = 0.0001, p = 0.02$ across our experiments.

## 4.4 Actuator Degradation Adaptation Transformer

Drawing inspiration from Yu et al. [45], we propose Actuator Degradation Adaptation Transformer (ADAPT). Similar to the embodiment term $e_i$ used in EAT [45], ADAPT incorporates the actuator degradation rate $d^i$ as an element of the robot's trajectory, serving as a priori information to aid decision-making. This symbolizing enables us to infuse prior information into our experiment, curtailing the exploration space, as discussed in Section 4.1.

We chose local position encoding, distinct from prior research employing a global position encoding [9, 22, 45]. To be more specific, in prior methodologies, the position information to be encoded can be denoted as $t$ with its value within $[0, T)$, where $T$ consistently represented the maximum experiment length. Such an approach aims to introduce global sequence information to address challenges in other tasks. However, we argue that position encoding should be limited to the input sequence's context length. Thus, in our ADAPT, we set $t \in (0, T_c)$, with $T_c$, the context length, typically much shorter than $T$, aligning more closely with the reality of robot locomotion. This design avoids potential timestep distribution anomalies, especially during prolonged walks, and enhances experimental performance.

We adopt an offline training approach similar to GATO [33] to train ADAPT, utilizing basic behavior cloning. The objective function is formally defined as follows:

$$\mathcal{L}(\theta) = \sum_{t=1}^{T} \|\hat{a}_t, \bar{a}_t\|_2^2, \tag{11}$$
$$(\hat{a}_1, \hat{a}_2, \dots, \hat{a}_T) \sim \text{Trans}_\theta(s_1, d_1, \bar{a}_1, s_2, d_2, \bar{a}_2, \dots, s_T, d_T, \bar{a}_T).$$

where $(s_t, d_t, \bar{a}_t)_{t=1,2,\dots T}$ is the sequence of transitions sampled in $\mathcal{T}$. While GATO uses MLE loss to characterize the distance to experts, we adopt MSE loss instead because of the continuous action space. The overall framework of ADAPT is shown in Figure 2 [1].

---

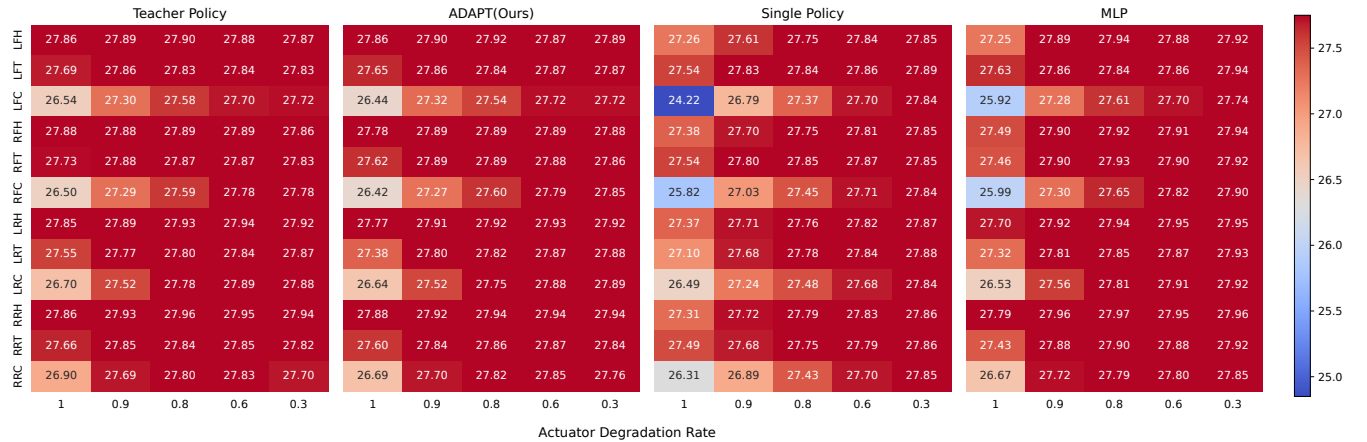[1] The source code has been open-sourced at https://github.com/WentDong/Adapt.

**Figure 3: Policy performance over different fault settings (joint, actuator degradation rate). The horizontal axis signifies the actuator degradation rate, with 1 indicating complete damage, while the vertical axis corresponds to the specific joint affected. Each grid in the figure represents the accumulative rewards averaged over 1024 runs in parallel simulation for that specific scenario. Each run randomly sampled the initial state, command speed, and fault occurrence time within specified ranges. The reward recording stopped if the robot fell or 1000 timesteps were collected.**

## 5 EXPERIMENT

In our cognition, few people are involved in this research topic. Therefore, our experiments focus on answering the following research questions (RQs):

**RQ1** How does the model perform under diverse actuator degradation scenarios?

**RQ2** What are the specific mechanisms by which the model can adapt to changes in actuator degradation?

**RQ3** What is the performance of our algorithm when deployed to a real-world robot?

For ease of presentation, we adhere to the joint naming conventions established by Unitree, which employ a three-letter scheme for each joint. The first letter of joint signifies **F**ront or **R**ear, the second letter of joint denotes **L**eft or **R**ight, and the final letter indicates **H**ip, **T**high or **C**alf. For example, the left front calf is abbreviated as LFC. More experiment settings can be seen in Appendix A

### 5.1 Evaluation in Simulator

To answer **RQ1**, we extensively evaluated our model's performance in the simulation environment. The results are shown in Figure 3. Due to the lack of existing strategies in this field, we established the following baselines for comparison (to ensure fairness, they were all trained in an identical environment setup):

- **Single Policy.** We tried to learn the same skills with one unified model, denoted as Single Policy, in our experiments. The Single Policy was trained via RL under the conditions corresponding to the sum of all 12 teacher policies until convergence.
- **MLP.** As an alternative approach, we replaced the transformer architecture in our ADAPT with an MLP, using the same training data, resulting in the model denoted as MLP.
- **Teacher Policy.** We also assessed the performance of the teacher policy under the same conditions, serving as an oracle to evaluate the learning capacity of our student policy.

We present a comparative analysis of the performance achieved by four distinct control policies within a simulated environment, visualized through a heat map in Figure 3. The control strategy corresponding to each figure is indicated at the top. We did not uniformly select points on the horizontal axis because, in our experiments, when the degradation rate is low, all models show subtle performance variation in response to changes in the degradation rate. Therefore, sparsely chosen points are sufficient to represent their performance trends. For more detailed information regarding how performance changes with degradation rate, please refer to Appendix B.

The comparison demonstrates ADAPT's adaptability and potential for zero-shot transfer into real-world scenarios marked by diverse joint actuator faults. When the actuator degradation rate is low (e.g., below 0.6), all the baselines can achieve satisfactory performance since the models can easily generalize such scenarios through domain randomization. However, when the actuator degradation rate rises to some extent, e.g., greater than 0.9, both Single Policy and MLP suffer a severe performance deterioration. Impressively, our model closely matches the Teacher Policy's performance across most scenarios and sometimes even outperforms them. Notably, our model's parameters are significantly less than 12 teacher policies, highlighting the pivotal role of knowledge distillation. Specifically, the total parameters of our model are 1.2M, while each teacher policy has 0.4M parameters, adding up 4.8M. To make our experiment convincing, we also record the standard deviation of the model's scores under each scenario in the experiment above. Our approach has the most stable score, and the detailed experimental results are provided in Appendix D.

Furthermore, the figure indicates that adapting to faults in the calf joints proves more difficult under identical training conditions (rows 3, 6, 9, 12 in Figure 3). This phenomenon is intuitive since the calf joints must exert larger torque to maintain stability and mobility. Therefore, they are significantly more affected by actuator degradation. We believe that these joints' adaptability to actuator
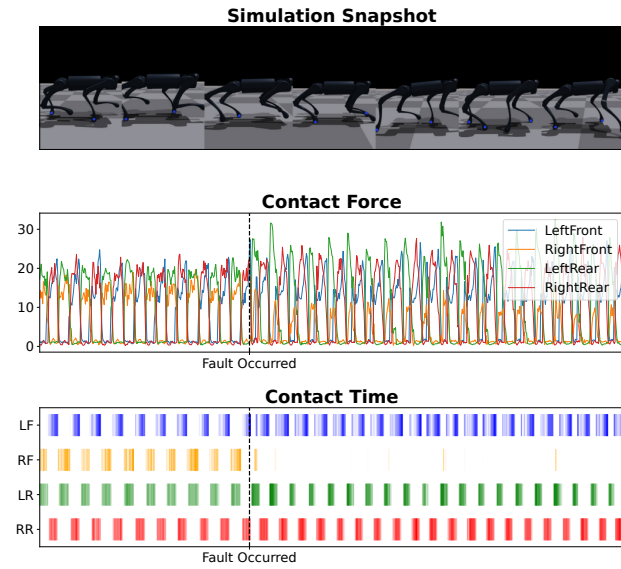
**Simulation Snapshot**



**Gravity Projection Over Time**



Figure 5: The projection of gravity components on the robot's plane during robot locomotion, with the robot's forward direction denoted as the positive y-axis and facing right as the positive x-axis. Different colors reflect changes in gravitational force components over time.

Figure 4: Illustration of Adaptive gait. Top: The robot's gait changes in the simulator, with ground-contact feet in blue dots. Middle: The L2-Norm of contact force on different feet. Bottom: feet contact time with the ground; darker color represents larger ground contact force.

degradation is more crucial. Consequently, their performance in severe actuator degradation scenarios in this experiment is worth noting, as it offers valuable insights into the effectiveness of the control policy.

Meanwhile, exploring the limits of our method's capabilities is also interesting. We tested the control performance of our ADAPT when multiple joints failed simultaneously, which means more severe damage or even multiple legs can not work properly at the same time. This situation is equivalent to a sudden drastic change in the external environment for the model. The experiment results show that our method can adapt to this kind of sudden change to a certain extent. Detailed experimental conclusions can be found in Appendix E.

## 5.2 Impaired Adaptive Gait Analysis

We find that our robot adapts to different joint actuator degradation scenarios through what we term "impaired adaptive gait", which describes our robot's capacity to modify its walking pattern in response to varying degrees of joint degradation or impairment. To address **RQ2**, we primarily investigate how our model accommodates regular and impaired adaptive gait patterns in this section.

Specifically, we command the robot to walk forward at a speed of 1 m/s and then abruptly induce complete degradation in one of its joints. We recorded the pressure feedback from its foot sensors and analyzed gait variations based on this force sensing, as depicted in Figure 4. To pose a challenge, we set the forward velocity to the maximum linear velocity limit in the training data and select the RFC joint for degradation, as evident from Figure 3. Notably, gait transitions may not be as pronounced in more straightforward tasks.
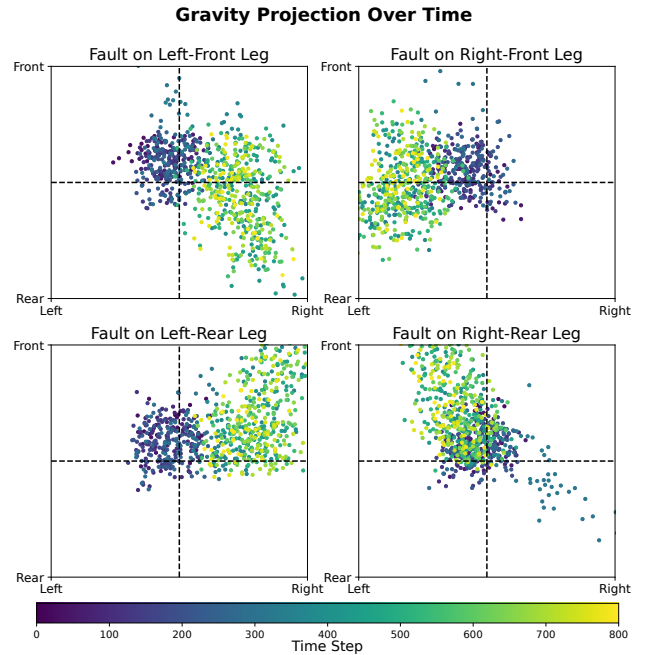
It can be observed that after the fault occurs, the contact force on the right front (RF) leg quickly decreases, which suggests that the right front leg is dragged forward and cannot provide support. In contrast, the left rear leg (LR) and right rear (RR) legs show a significant darkening in color, indicating that they provide more support (which is also evident in the middle image). At the same time, the left front (LF) leg has a longer contact time, which means the robot's center of gravity shifts, allowing it to achieve a new stable state under the support of three legs and continue its walking task.

Interestingly, we discovered that when gait transitions occur, the robot's center of gravity shifts towards the opposite side of the faulty joint, which is against the intuition that the center of gravity would move towards the direction of the faulty joint due to a lack of support in the direction. To verify this, we conducted experiments in a noisy environment, where the command speed was randomly sampled between 0.2 and 1.0, and the faults occurred randomly between 200 and 500 timesteps. We plotted the gravity projection at different timesteps in Figure 5.

It can be observed that at the beginning of each episode, the gravitational force components are initially close to the origin. As time progresses and faults occur, their gravity projection gradually shifts toward a particular direction. The four plots correspond to severe calf joint actuator degradation scenarios for each leg. It can be observed that when a fault occurs, the direction of the center of gravity shift tends to be away from the location of the fault leg. In other words, the model adjusts its gait to shift the center of gravity toward a position that facilitates three-legged support. This phenomenon is evident across various joint degradation scenarios and
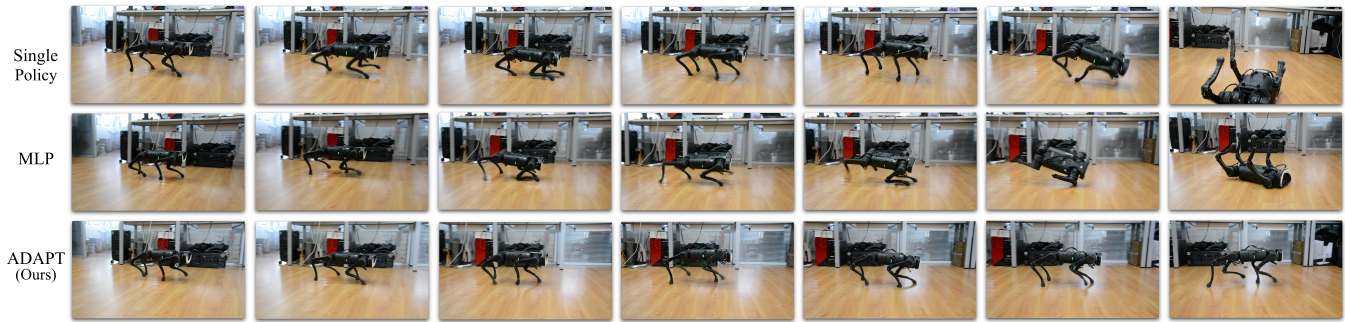
**Figure 6: Real-world deployment results of three control models for actuator degradation, with damage restricted to the LFC joint. Images have been mirrored for visual consistency.**

becomes more pronounced as the degree of degradation increases. For additional experimental results, please refer to Appendix C. This is a strong indicator that our control strategy possesses a distinct impaired adaptive gait

## 5.3 Real World Deployment

This section primarily delves into addressing **RQ3**. Our work can be successfully transferred to a real-world quadruped robot, demonstrating high adaptability to actuator degradation that aligns with its simulation performance.[2]

To further demonstrate the effectiveness of our proposed model, we directly deploy the models of Single Policy, MLP, and our "Adapt" on a Unitree A1 robot in the real world and induce multiple degradation rates across all joints. We simulate actual actuator degradation by attenuating the model's current output torques before sending them to the actuators. Though successful in low degradation rate scenarios, both Single Policy and MLP fail when the degradation rate is greater than or equal to 0.9.

Figure 6 illustrates the performance of three different models when adapting to the left front calf (LFC) joint experiencing a degradation rate of 0.9. [3] As the figure shows, the Single Policy model exhibits no signs of gait transition. Upon degradation, it attempts to drive the faulty leg similarly to the original gait to maintain balance, leading to imbalance and a fall. Conversely, MLP begins the transition to an adaptive gait only after experiencing a significant deviation in the base to maintain balance. As a result, it still encounters difficulties in maintaining stability. In stark contrast, our "Adapt" method rapidly adjusts the gait by shifting the center of mass to the opposite side (consistent with our analysis in Section 5.2) and completes the walking task.

According to real-world experiments, the Single Policy model does not learn the adaptability required for severe actuation degradation, even though it underwent the same training scenarios. The MLP baseline lacks the sensitivity to perceive sudden state changes due to the absence of historical information. It only transfers to an adaptive gait when there is a significant deviation from its normal state, which is too late to regain balance. Our Adapt shows the capability for zero-shot transfer to real-world robot actuator degradation scenarios successfully.

## 6 CONCLUSION AND FUTURE WORK

This paper introduces the Adapt framework to enhance quadruped robots' adaptability to joint actuator faults. Our approach utilizes reinforcement learning and a well-designed architecture for smooth transitions between intact and faulty states. It can intelligently adapt to various impaired adaptive gaits when an actuator degradation fault occurs. Notably, we demonstrate successful zero-shot transferability from simulator to real-world robots, showcasing adaptability across different scenarios. In addition, our work has demonstrated the ability to resist sudden environmental changes caused by multiple joint failures to a certain extent. As robotics progresses, our framework could further contribute to robust locomotion and broader adaptability across industries.

Despite this, there is still much room for improvement in our work. On the one hand, our approach focuses on a specific quadruped robot model, but extending it to various robotic platforms can further demonstrate its generality. On the other hand, this work focuses more on the motion control part of the quadruped robot, integrating fault detection methods into the control model for a more unified solution is a promising direction.

## REFERENCES

[1] Maxime Allard, Simón C Smith, Konstantinos Chatzilygeroudis, Bryan Lim, and Antoine Cully. 2023. Online Damage Recovery for Physical Robots with Hierarchical Quality-Diversity. *ACM Transactions on Evolutionary Learning* 3, 2 (2023), 1–23.
[2] Kai Arulkumaran, Marc Peter Deisenroth, Miles Brundage, and Anil Anthony Bharath. 2017. Deep reinforcement learning: A brief survey. *IEEE Signal Processing Magazine* 34, 6 (2017), 26–38.
[3] C Dario Bellicoso, Marko Bjelonic, Lorenz Wellhausen, Kai Holtmann, Fabian Günther, Marco Tranzatto, Peter Fankhauser, and Marco Hutter. 2018. Advances in real-world applications for legged robots. *Journal of Field Robotics* 35, 8 (2018), 1311–1326.
[4] Matteo Bettini, Ajay Shankar, and Amanda Prorok. 2023. Heterogeneous Multi-Robot Reinforcement Learning. arXiv:2301.07137 [cs.RO]
[5] Priyaranjan Biswal and Prases K Mohanty. 2021. Development of quadruped walking robots: A review. *Ain Shams Engineering Journal* 12, 2 (2021), 2017–2031.
[6] Adrian Boeing and Thomas Bräunl. 2012. Leveraging multiple simulators for crossing the reality gap. In *2012 12th international conference on control automation*

---

[2]More real-world videos can be found at https://sites.google.com/view/adapt-2023.
[3]We did not directly set the degradation rate to 1 because, in such case, the joint actuator would show noticeable damping due to the ongoing motor. However, this does not accurately represent real-world conditions with actuator degradation.

robotics & vision (ICARCV). IEEE, 1113–1119.

[7] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. 2016. Openai gym. arXiv preprint arXiv:1606.01540 (2016).

[8] Jinlin Chen, Jiannong Cao, Zhiqin Cheng, and Wei Li. 2023. Mitigating Imminent Collision for Multi-Robot Navigation: A TTC-Force Reward Shaping Approach. In Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems (London, United Kingdom) (AAMAS '23). International Foundation for Autonomous Agents and Multiagent Systems, Richland, SC, 1448–1456.

[9] Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. 2021. Decision transformer: Reinforcement learning via sequence modeling. Advances in neural information processing systems 34 (2021), 15084–15097.

[10] Zhijun Chen, Qingxing Xi, Feng Gao, and Yue Zhao. 2022. Fault-tolerant gait design for quadruped robots with one locked leg using the GF set theory. Mechanism and Machine Theory 178 (2022), 105069.

[11] Junwen Cui, Zhan Li, Jing Qiu, and Tianxiao Li. 2022. Fault-tolerant motion planning and generation of quadruped robots synthesised by posture optimization and whole body control. Complex & Intelligent Systems 8, 4 (2022), 2991–3003.

[12] Jared Di Carlo, Patrick M Wensing, Benjamin Katz, Gerardo Bledt, and Sangbae Kim. 2018. Dynamic locomotion in the mit cheetah 3 through convex model-predictive control. In 2018 IEEE/RSJ international conference on intelligent robots and systems (IROS). IEEE, 1–9.

[13] Alejandro Escontrela, Xue Bin Peng, Wenhao Yu, Tingnan Zhang, Atil Iscen, Ken Goldberg, and Pieter Abbeel. 2022. Adversarial motion priors make good substitutes for complex reward functions. In 2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). IEEE, 25–32.

[14] Junsong FAN, Yuxi WANG, He GUAN, Chunfeng SONG, and Zhaoxiang ZHANG. 2022. Toward few-shot domain adaptation with perturbation-invariant representation and transferable prototypes. Frontiers of Computer Science 16, 3, Article 163347 (2022). https://doi.org/10.1007/s11704-022-2015-7

[15] Dongdong Guo, Li Fu, and Lingling Wang. 2019. Robots solving the urgent problems by themselves: A review. In 2019 IEEE International Conference on Consumer Electronics-Taiwan (ICCE-TW). IEEE, 1–2.

[16] Jonathan Ho and Stefano Ermon. 2016. Generative adversarial imitation learning. Advances in neural information processing systems 29 (2016).

[17] Carolin Kemper and Michael Kolain. 2022. K9 Police Robots-Strolling Drones, RoboDogs, or Lethal Weapons?. In Accepted paper at WeRobot2022 conference.

[18] Eliahu Khalastchi and Meir Kalech. 2018. A sensor-based approach for fault detection and diagnosis for robotic systems. Autonomous Robots 42 (2018), 1231–1248.

[19] Sylvain Koos, Jean-Baptiste Mouret, and Stéphane Doncieux. 2010. Crossing the reality gap in evolutionary robotics by promoting transferable controllers. In Proceedings of the 12th annual conference on Genetic and evolutionary computation. 119–126.

[20] Eric Krotkov, Douglas Hackett, Larry Jackel, Michael Perschbacher, James Pippine, Jesse Strauss, Gill Pratt, and Christopher Orlowski. 2017. The DARPA Robotics Challenge Finals: Results and Perspectives. Journal of Field Robotics 34, 2 (2017), 229–240. https://doi.org/10.1002/rob.21683 arXiv:https://onlinelibrary.wiley.com/doi/pdf/10.1002/rob.21683

[21] Ashish Kumar, Zipeng Fu, Deepak Pathak, and Jitendra Malik. 2021. Rma: Rapid motor adaptation for legged robots. arXiv preprint arXiv:2107.04034 (2021).

[22] Hang Lai, Weinan Zhang, Xialin He, Chen Yu, Zheng Tian, Yong Yu, and Jun Wang. 2023. Sim-to-Real Transfer for Quadrupedal Locomotion via Terrain Transformer. In 2023 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 5141–5147.

[23] Joonho Lee, Jemin Hwangbo, Lorenz Wellhausen, Vladlen Koltun, and Marco Hutter. 2020. Learning quadrupedal locomotion over challenging terrain. Science robotics 5, 47 (2020), eabc5986.

[24] Dikai Liu, Tianwei Zhang, Jianxiong Yin, and Simon See. 2022. Saving the Limping: Fault-tolerant Quadruped Locomotion via Reinforcement Learning. arXiv preprint arXiv:2210.00474 (2022).

[25] Viktor Makoviychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles Macklin, David Hoeller, Nikita Rudin, Arthur Allshire, Ankur Handa, et al. 2021. Isaac gym: High performance gpu-based physics simulation for robot learning. arXiv preprint arXiv:2108.10470 (2021).

[26] Gabriel B Margolis, Ge Yang, Kartik Paigwar, Tao Chen, and Pulkit Agrawal. 2022. Rapid locomotion via reinforcement learning. arXiv preprint arXiv:2205.02824 (2022).

[27] Takahiro Miki, Joonho Lee, Jemin Hwangbo, Lorenz Wellhausen, Vladlen Koltun, and Marco Hutter. 2022. Learning robust perceptive locomotion for quadrupedal robots in the wild. Science Robotics 7, 62 (2022), eabk2822.

[28] I Made Aswin Nahrendra, Byeongho Yu, and Hyun Myung. 2023. DreamWaQ: Learning robust quadrupedal locomotion with implicit terrain imagination via deep reinforcement learning. In 2023 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 5078–5084.

[29] Wataru Okamoto and Kazuhiko Kawamoto. 2020. Reinforcement learning with randomized physical parameters for fault-tolerant robots. In 2020 Joint 11th

[30] Xue Bin Peng, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. 2018. Sim-to-real transfer of robotic control with dynamics randomization. In 2018 IEEE international conference on robotics and automation (ICRA). IEEE, 3803–3810.

[31] Xue Bin Peng, Erwin Coumans, Tingnan Zhang, Tsang-Wei Lee, Jie Tan, and Sergey Levine. 2020. Learning agile robotic locomotion skills by imitating animals. arXiv preprint arXiv:2004.00784 (2020).

[32] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. Improving language understanding by generative pre-training. (2018).

[33] Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gomez Colmenarejo, Alexander Novikov, Gabriel Barth-Maron, Mai Gimenez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, et al. 2022. A generalist agent. arXiv preprint arXiv:2205.06175 (2022).

[34] Unitree Robotics. 2020. Unitree A1. Online. https://www.unitree.com/a1

[35] Alejandro Romero, Gianluca Baldassarre, Richard J. Duro, and Vieri Giuliano Santucci. 2023. Learning Multiple Tasks with Non-Stationary Interdependencies in Autonomous Robots. In Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems (London, United Kingdom) (AAMAS '23). International Foundation for Autonomous Agents and Multiagent Systems, Richland, SC, 2547–2549.

[36] Nikita Rudin, David Hoeller, Philipp Reist, and Marco Hutter. 2022. Learning to Walk in Minutes Using Massively Parallel Deep Reinforcement Learning. In Proceedings of the 5th Conference on Robot Learning (Proceedings of Machine Learning Research, Vol. 164), Aleksandra Faust, David Hsu, and Gerhard Neumann (Eds.). PMLR, 91–100. https://proceedings.mlr.press/v164/rudin22a.html

[37] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347 (2017).

[38] Prasanth Sengadu Suresh, Yikang Gui, and Prashant Doshi. 2023. Dec-AIRL: Decentralized Adversarial IRL for Human-Robot Teaming. In Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems (London, United Kingdom) (AAMAS '23). International Foundation for Autonomous Agents and Multiagent Systems, Richland, SC, 1116–1124.

[39] Jie Tan, Tingnan Zhang, Erwin Coumans, Atil Iscen, Yunfei Bai, Danijar Hafner, Steven Bohez, and Vincent Vanhoucke. 2018. Sim-to-real: Learning agile locomotion for quadruped robots. arXiv preprint arXiv:1804.10332 (2018).

[40] Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel. 2017. Domain randomization for transferring deep neural networks from simulation to the real world. In 2017 IEEE/RSJ international conference on intelligent robots and systems (IROS). IEEE, 23–30.

[41] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. Advances in neural information processing systems 30 (2017).

[42] Muning WEN, Runji LIN, Hanjing WANG, Yaodong YANG, Ying WEN, Luo MAI, Jun WANG, Haifeng ZHANG, and Weinan ZHANG. 2023. Large sequence models for sequential decision-making: a survey. Frontiers of Computer Science 17, 6, Article 176349 (2023). https://doi.org/10.1007/s11704-023-2689-5

[43] Jinze Wu, Guiyang Xin, Chenkun Qi, and Yufei Xue. 2023. Learning Robust and Agile Legged Locomotion Using Adversarial Motion Priors. IEEE Robotics and Automation Letters (2023).

[44] Zichen Yan, Junbo Tan, Bin Liang, Houde Liu, and Jun Yang. 2022. Active Fault-Tolerant Control Integrated with Reinforcement Learning Application to Robotic Manipulator. In 2022 American Control Conference (ACC). IEEE, 2656–2662.

[45] Chen Yu, Weinan Zhang, Hang Lai, Zheng Tian, Laurent Kneip, and Jun Wang. 2023. Multi-embodiment Legged Robot Control as a Sequence Modeling Problem. In 2023 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 7250–7257.

[46] Wenhao Yu, Jie Tan, C Karen Liu, and Greg Turk. 2017. Preparing for the unknown: Learning a universal policy with online system identification. arXiv preprint arXiv:1702.02453 (2017).

[47] Yu Zhang. 2023. From Abstractions to Grounded Languages for Robust Coordination of Task Planning Robots. In Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems (London, United Kingdom) (AAMAS '23). International Foundation for Autonomous Agents and Multiagent Systems, Richland, SC, 2535–2537.

[48] Rui Zhao, Xu Liu, Yizheng Zhang, Minghao Li, Cheng Zhou, Shuai Li, and Lei Han. 2023. CraftEnv: A Flexible Collective Robotic Construction Environment for Multi-Agent Reinforcement Learning. In Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems (London, United Kingdom) (AAMAS '23). International Foundation for Autonomous Agents and Multiagent Systems, Richland, SC, 1164–1172.

[49] Yong Yong Zhao, Jing Hua Wang, Bao Wen Zhang, Xu Yao, and Guo Hua Cao. 2023. Gait planning and fault-tolerant control of quadruped robots. In Third International Conference on Mechanical, Electronics, and Electrical and Automation Control (METMS 2023), Vol. 12722. SPIE, 841–846.

International Conference on Soft Computing and Intelligent Systems and 21st International Symposium on Advanced Intelligent Systems (SCIS-ISIS). IEEE, 1–4.

# A TRAINING DETAILS

## A.1 Experiment Settings

**Simulation.** We implement our model and baselines based on the open-source codebase provided by Rudin et al. [36], which ensures thousands of robots simulated in parallel in the environment of IsaacGym simulator [25]. The URDF file of the A1 robot from Unitree is used during the simulation.

**Hardware.** We use the A1 robot from Unitree for all the real-world application experiments, which has 12 actuators (3 motors on each leg) with a maximum torque of 33.5N.m [34].



**Figure 7: The relationship between the performance of the four models and the degradation rate. The horizontal axis represents the actuator degradation rate. The vertical axis shows the model's performance.**

## A.2 Training Details

For the teacher policy, we employ a three-layer MLP (multi-layer perceptron) with $(512, 256, 128)$ hidden units to encode the states along with the degradation rates and produce the actions. We apply curriculum learning for teacher policy training, which gradually increases the range of the command velocity. Given the physical structure of the A1 robot, it is challenging for the quadruped robot to move laterally or backward with actuator degradation. Hence, we concentrate on forward walking with turning commands sufficient for omnidirectional locomotion [26]. The maximum walking velocity is set to 1.0 m/s, and the maximum turning angular velocity is 0.5 rad/s.

Moreover, the degradation rate distribution significantly affects the policy's performance. If we sample the degradation rate from a uniform distribution, the training process will be unstable, leading to poor performance under high degradation. Therefore, we use the adaptive sampling mechanism described in Section 4.3 for the degradation rate both in the teacher training and data collection phase, which ensures that the degradation changes according to rules that maintain a reasonable distribution.

We train each teacher policy on one NVIDIA GeForce RTX 3080 for 3000 iterations, which simultaneously takes about 75 minutes, parallelizing 4096 environments. Afterward, we adopt our proposed ADAPT framework for student policy training by setting context length $T_c = 20$ and the number of blocks $n_b = 6$. We train the ADAPT model on the collected dataset described in Section 4.3, which takes 500 million updates and 20 hours on a single NVIDIA Geforce RTX 3080 GPU.

## A.3 Domain Randomization

Domain Randomization [40] was proposed to address the simulator-reality gap, which randomizes the simulator parameters to expose the model to a wide range of environments during the training process. We adopt Domain Randomization during the training process of teacher policies as well as the collection process of the data set. The details of randomization can be viewed in Table 2.

**Table 2: Domain Randomization Range**

| Parameter | Range |
|---|---|
| Friction | [0.05, 2.75] |
| Restitution | [0, 1] |
| Base Mass | [-1, 2] |
| Link Mass | [0.8, 1.2] |
| COM Position | [-0.05, 0.05] |
| $K_p$ | [22.4, 33.6] |
| $K_d$ | [0.56, 0.84] |
| Motor Strength | [0.8, 1.2] |

# B PERFORMANCE VARIATION WITH DEGRADATION RATE

We conducted a series of experiments identical to Section 5.1, but with a much denser range of degradation rate variations, as shown in Figure 7.

Each subplot illustrates the performance variation of a model under different levels of degradation for 12 joints. The performance is measured by the average reward obtained during experiments. It's evident that all models exhibit relatively subtle performance changes when the degradation rate is below 0.8. However, when it exceeds 0.8, a significant performance decline becomes noticeable. This observation further corroborates the discussions in the earlier sections regarding task difficulty. Additionally, it's also evident that our ADAPT closely approaches the performance of the Teacher Policy when the degradation rate is low. However, at higher degradation rates, it notably outperforms the baseline. This observation is consistent with the experimental findings presented in Section 5.1.

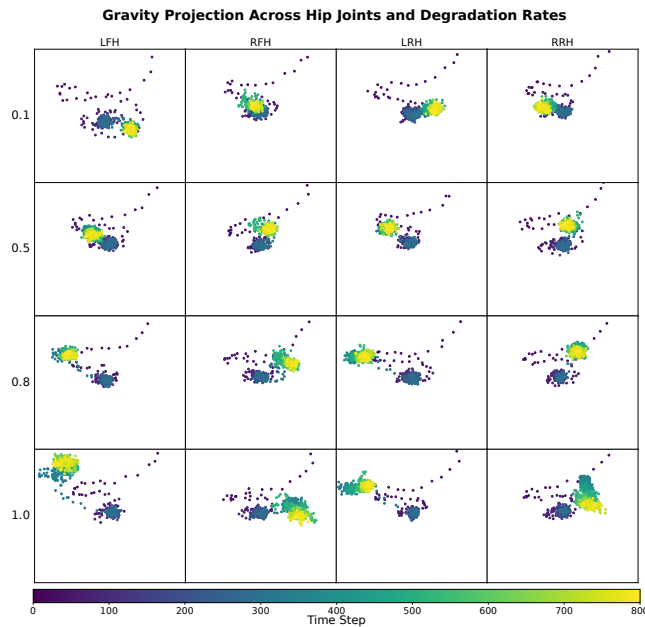# C IMPAIRED ADAPTIVE GAIT ANALYSIS OF VARIOUS SCENARIOS



**Figure 8: Each column represents a specific damaged joint within the same leg, and each row corresponds to the same degradation rate.**

To comprehensively discuss the gait variations of our control models under different degradation scenarios, we conducted experiments to monitor the shifts in the gravitational projection direction for all joint degradation scenarios. The experimental setup and presentation method are consistent with those in Figure 5. To visually emphasize changes in the gravitational projection under different degradation levels, we selected cases with degradation rates of 0.1, 0.5, 0.8, and 1.0 for presentation. This choice is based on discussions in Section 5.1 and Appendix B.

Figure 8 illustrates the outcomes of hip joint degradation scenarios. From the figures, it is evident that when the hip joint is degraded, as the level of degradation increases, the gravity projection tends to shift more towards the degraded side. Seldom reverse shift in the projection of gravity is observed, signifying that the hip joint does not transition to the impaired adaptive gait. This is because the hip joint primarily handles leg adduction and abduction and is not as critical for stabilizing the body as other joints. Its functionality can naturally be compensated for by other joints. Consequently, it doesn't require a transition into the impaired adaptive gait. Such an outcome reflects that the hip joint adaptation is less challenging in the presence of actuator faults, consistent with the conclusions in Figures 3 and 7.

Figure 9 and Figure 10 depict thigh and calf joint degradation scenarios, respectively. It's evident that as the degree of degradation increases, both scenarios exhibit a more significant shift in the gravity projection away from the faulty side. This suggests that they both show varying degrees of impaired adaptive gait. Therefore, the points in Figure 10 appear smaller because, compared to the thigh joint, the shift of gravity projection is more significant when
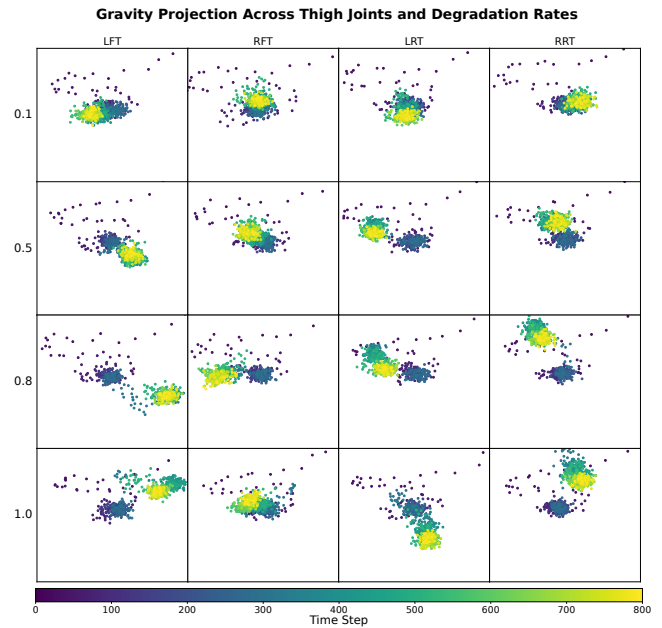


**Figure 9: Each column represents a specific damaged joint within the same leg, and each row corresponds to the same degradation rate.**

the calf joint is degraded. This reflects that the calf joint undergoes more noticeable gait changes. This observation aligns with the earlier conclusion drawn from our experiments, emphasizing the more significant challenge of adapting the calf joint to degradation scenarios.
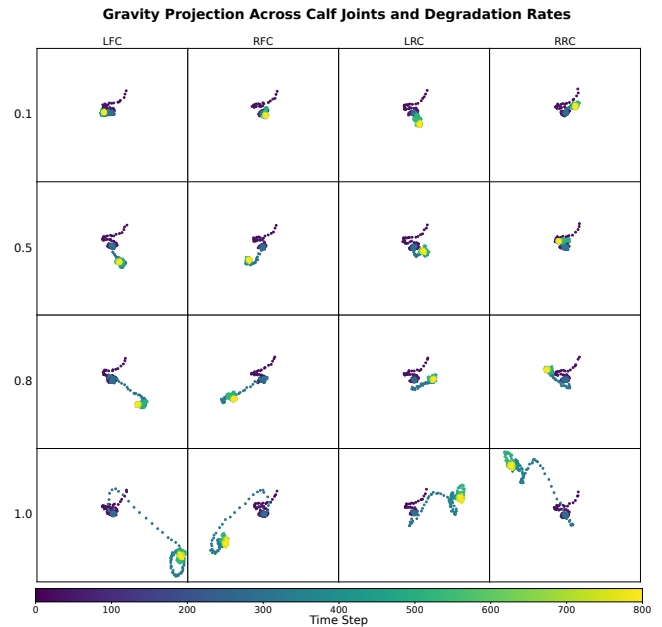


**Figure 10: Each column represents a specific damaged joint within the same leg, and each row corresponds to the same degradation rate.**
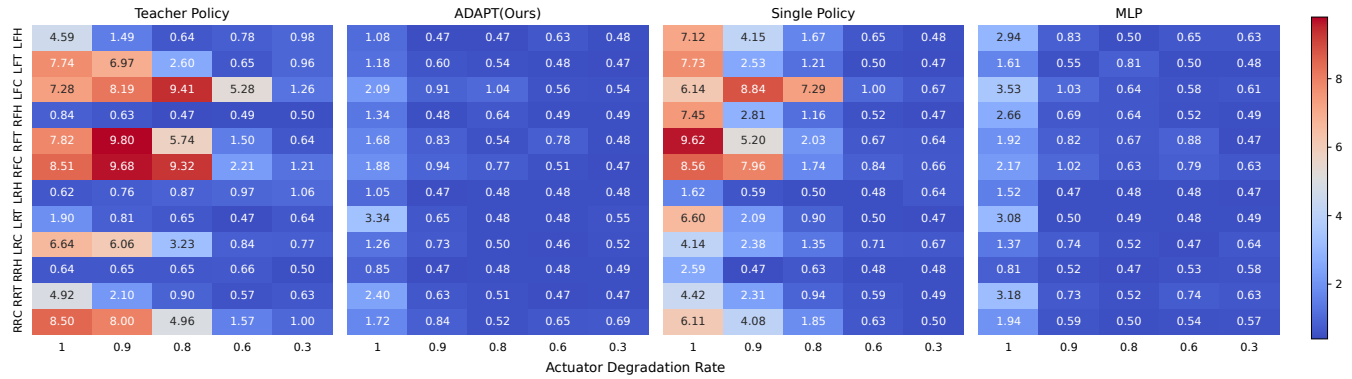
**Figure 11: The horizontal axis signifies the actuator degradation rate, while the vertical axis corresponds to the specific joint affected. Each grid in the figure represents the standard deviation of rewards over 1024 runs in parallel simulation for that specific scenario.**

## D COMPARISON OF STANDARD DEVIATION OF SCORES UNDER DIFFERENT SCENARIOS

We experimented with the same settings as in Figure 3 and recorded the standard deviation of scores in all conditions. The result is shown in Figure 11. From the figure, it can be seen that as the actuator degradation rate increases, the standard deviation of the model's score becomes larger, indicating that the model's performance is more unstable as it faces more difficult problems. However, the Teacher Policy model performs significantly unstable, which we believe is because it has seen relatively few fault situations during training, which limits its generalization ability. The Single Policy also exhibits a relatively large standard deviation because it cannot effectively handle fault situations, which is consistent with other experiments' conclusions. Both MLP and our ADAPT show relatively small standard deviations, which means that our models perform consistently better than MLP in critical joints.

## E THE SCENARIO OF MULTIPLE JOINTS GETTING INTO FAULT SIMULTANEOUSLY

Due to the completely different state space changes caused by multi-joint faults compared to single-joint faults, it can be considered as a completely new task. It should be noted that the above fault scenarios are rarely encountered in real-world deployments. These experiments aim to demonstrate that our method can directly generalize to adapt to certain multi-joint fault scenarios without fine-tuning.

### E.1 Dual-Joint Actuator Degradation

Figure 12 represents the performance of a robot when encountering simultaneous failure of two joints measured by the average reward obtained from 1024 parallel experiments in the simulator. This figure depicts the robot's performance when encountering simultaneous faults in two joints. As can be seen from the figure, even when 90% of any two joints are damaged, our model can still achieve good average rewards in more than 90% of the cases.

### E.2 Three-Joint Actuator Degradation

We recorded the average reward and the time taken for the robot to maintain its balance after any three joints occurred actuator
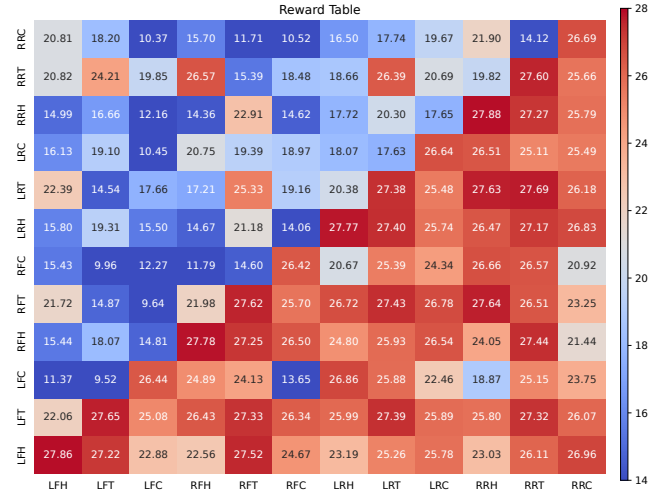


**Figure 12: Each grid represents the average reward, with the horizontal and vertical axes representing the specific fault joints. The upper left area represents the case where both joints' actuator degradation rate is 1, while the lower right area represents the case where both joints have an actuator degradation rate of 0.9.**

degradation simultaneously in 512 experiments. We plotted the distribution in Figure 13. It can be seen that when the actuator degradation rate reaches 70%, the average reward performance of the model remains at a high level. When the actuator degradation rate reaches 80%, the average reward decreases, but the survival time does not change significantly, with only an average 6.5% decrease compared to the original performance. This means that in this situation, the model can still maintain its stability but may not be able to strictly follow commands.

### E.3 All Joint Actuator Degradation

Finally, we tested our model's reward and the time taken to maintain balance after all joints simultaneously experienced a certain degree of malfunction. The data was obtained by averaging results from 1024 parallel experiments and shown in Figure 14. It can be seen that when the degradation level of all actuators is below 25%, the model
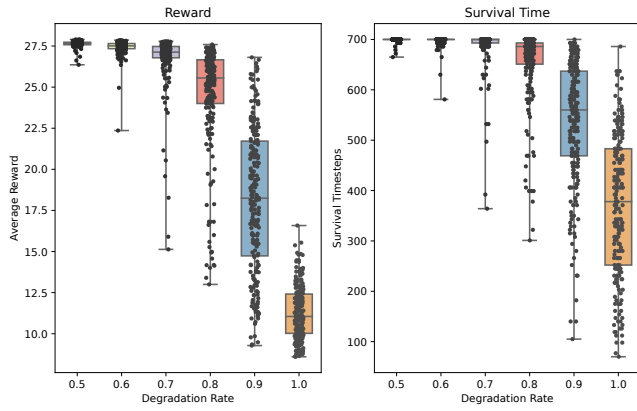
**Figure 13: The left graph represents the average reward obtained from 512 parallel environments, and the right graph represents the time taken by the robot to maintain its stability after encountering actuator degradation.**

performance is not significantly affected, with a decrease of less than 2.5% in performance compared to the original level. However, as the actuator degradation rate increases, the robot's survival time fluctuates, which is because when the actuator degradation rate increases to a certain level, the robot gives up its movement ability and prioritizes maintaining its standing state.

Let us consider actuator health as an external environmental indicator. Our model performance will not be significantly affected as long as the change in this indicator is less than or equal to 25%. However, further exploration is needed to determine the specific boundaries of the model's capabilities.
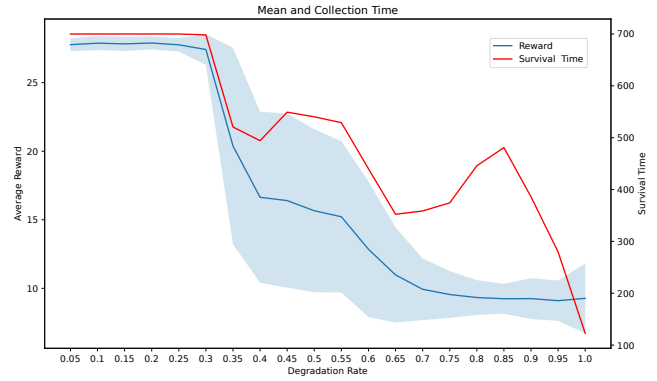


**Figure 14: The horizontal axis represents the actuator degradation rate, with the blue line representing the average reward and the blue shadow reflecting their variance. The red line represents the time taken for the robot to maintain its balance after the failure occurs.**